

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

Meanders as a scaling motif for understanding of floodplain soil microbiome and biogeochemical potential at the watershed scale

Paula B. Matheus Carnevali¹, Adi Lavy¹, Alex D. Thomas², Alexander Crits-Christoph³, Spencer Diamond¹, Raphaeël Meéheust^{1,4}, Matthew R. Olm^{3,^}, Allison Sharrar¹, Shufei Lei¹, Wenming Dong⁵, Nicola Falco⁵, Nicholas Bouskill⁵, Michelle Newcomer⁵, Peter Nico⁵, Haruko Wainwright⁵, Dipankar Dwivedi⁵, Kenneth H. Williams⁵, Susan Hubbard⁵, Jillian F. Banfield^{1,2,3,4,5,6,*}.

¹Department of Earth and Planetary Science, University of California, Berkeley, CA, USA.

²Department of Environmental Science, Policy, and Management, University of California, Berkeley, CA, USA.

³Department of Plant and Microbial Biology, University of California, Berkeley, CA, USA.

⁴Innovative Genomics Institute, Berkeley, CA, USA.

⁵Earth and Environmental Sciences Area, Lawrence Berkeley National Laboratory, Berkeley, CA, USA

⁶Chan Zuckerberg Biohub, San Francisco, CA, USA.

Current affiliation:

[^]Department of Microbiology and Immunology, Stanford University, Palo Alto, CA, USA.

*Corresponding author: jbanfield@berkeley.edu

42 **Abstract**

43 Biogeochemical exports of C, N, S and H₂ from watersheds are modulated by the activity of
44 microorganisms that function over micron scales. This disparity of scales presents a substantial
45 challenge for development of predictive models describing watershed function. Here, we tested
46 the hypothesis that meander-bound regions exhibit patterns of microbial metabolic potential that
47 are broadly predictive of biogeochemical processes in floodplain soils along a river corridor. We
48 intensively sampled floodplain soils located in the upper, middle, and lower reaches of the East
49 River in Colorado and reconstructed 248 draft quality genomes representative at a sub-species
50 level. Approximately one third of the representative genomes were detected across all three
51 locations with similar levels of abundance, and despite the very high microbial diversity and
52 complexity of the soils, ~15% of species were detected in two consecutive years. A core floodplain
53 microbiome was enriched in bacterial capacities for aerobic respiration, aerobic CO oxidation, and
54 thiosulfate oxidation with the formation of elemental sulfur. We did not detect systematic patterns
55 of gene abundance based on sampling position relative to the river. However, at the watershed
56 scale meander-bound floodplains appear to serve as scaling motifs that predict aggregate capacities
57 for biogeochemical transformations in floodplain soils. Given this, we conducted a transcriptomic
58 analysis of the middle site. Overall, the most highly transcribed genes were *amoCAB* and *nxrAB*
59 (for nitrification) followed by genes involved in methanol and formate oxidation, and nitrogen and
60 CO₂ fixation. Low soil organic carbon correlated with high activity of genes involved in methanol,
61 formate, sulfide, hydrogen, and ammonia oxidation, nitrite oxidoreduction, and nitrate and nitrite
62 reduction. Thus, widely represented genetic capacities did not predict *in situ* activity at one time
63 point, but rather they define a reservoir of biogeochemical potential available as conditions change.
64

65

66 **Introduction**

67 Watersheds are geographic areas that capture precipitation that is ultimately discharged
68 into rivers and other larger water bodies. Of particular interest are watersheds in mountainous
69 regions, as these are major sources of freshwater^{1,2}. Within mountainous watersheds, complex
70 interactions among vegetation, hydrology, geochemistry, and geology occur within and across
71 watershed compartments, including across bedrock-soil-vegetation compartments of terrestrial
72 hillslopes, across terrestrial-aquatic interfaces and within the fluvial system itself. Interactions
73 within a reactive watershed typically vary as a function of disturbance as well as landscape position
74 and topography. For example, interactions in an alpine region of a mountainous watershed are
75 likely to be quite different from a lower montane floodplain region³. Floodplains, which extend
76 from the river banks to the base of hillslopes, comprise the riparian zone (a vegetated interface
77 between the river channel and the rest of the ecosystem), and are notable as they integrate inputs
78 from all watershed compartments. They also display depositional gradients and features associated
79 with past and current river channel positions. Unlike hillslopes, floodplains receive water and
80 constituents either by surface runoff or groundwater discharge. They are typically significantly

81 impacted by changes in river conditions and can be inundated when river flow and stage increases
82 following snowmelt. Consequently, floodplains are dynamic compartments in which
83 hydrobiogeochemical processes vary seasonally and potentially spatially. Overall, floodplains are
84 important watershed regions in which microbial activity can modulate the form and abundance of
85 nutrients and contaminants derived from hillslopes and river water prior to their export from the
86 watershed.

87 Here, we conducted a study of floodplain soils of the mountainous East River (CO)
88 watershed to investigate how patterns in the distribution of soil microorganisms and their
89 associated functions and activities can induce geochemical gradients that impact riverine nutrient
90 and contaminant fluxes. We tested a ‘system-of systems’ approach⁴ wherein meander-bound
91 regions were selected as scaling motifs (repeating patterns along the river that can be used for
92 ecosystem modeling at the watershed scale), in which microbially-mediated biogeochemical
93 processes that are shaped by reactions occurring at the micron-scale might be representative of
94 processes throughout the floodplain. Detailed analyses of meander-bound floodplain soils may
95 reveal patterns that approximate watershed processes at the tens of kilometers scale, and could
96 provide much needed input for watershed hydrobiogeochemical models. This study applied
97 genome-resolved metagenomic and metatranscriptomic bioinformatics methods to large nucleic
98 acid sequence datasets to investigate microbial community composition and distribution and to
99 infer capacities for microbially-mediated C, S, H and N cycling and *in situ* activity in floodplain
100 soil microbial communities.

101

102 **Results**

103

104 *Metagenomes overview*

105 Three meander-bound floodplains following the meandering pattern of the East River (**Fig.**
106 **1a**) were chosen for this study: one upstream (meander-bound floodplain G (Floodplain G); **Fig.**
107 **1b**), one midstream (meander-bound floodplain L (Floodplain L); **Fig. 1c**), and one downstream
108 (meander-bound floodplain Z (Floodplain Z); **Fig. 1d**). Sample number was prioritized over
109 sequencing depth to better resolve the types and distribution patterns of the most abundant
110 organisms across the meander-bound floodplains (‘floodplains’ subsequently). An average 6.4
111 giga base pairs (Gbp; 3.2 – 11.5 Gbp) of sequencing data was obtained from 90 DNA extractions
112 out of 94 floodplain soil samples collected in 2015. An average 12 Gbp (6.0 – 15.0 Gbp) of

113 sequencing data was obtained from the other four samples. In total, ~0.6 Tbp of DNA sequence
114 information was acquired from samples collected in 2015. Our strategy aimed to capture the most
115 abundant members of the microbial community (instead of the whole community), so it is not
116 surprising that an average 13% (3 - 30%) of reads mapped to their respective assemblies
117 ([Supplementary Table 1](#)). This result also reflects the large (but variable) tail on the abundance
118 distribution of microorganisms in soil ⁵; and the communities captured by the smaller assemblies
119 resulting from samples with lower sequencing depth.

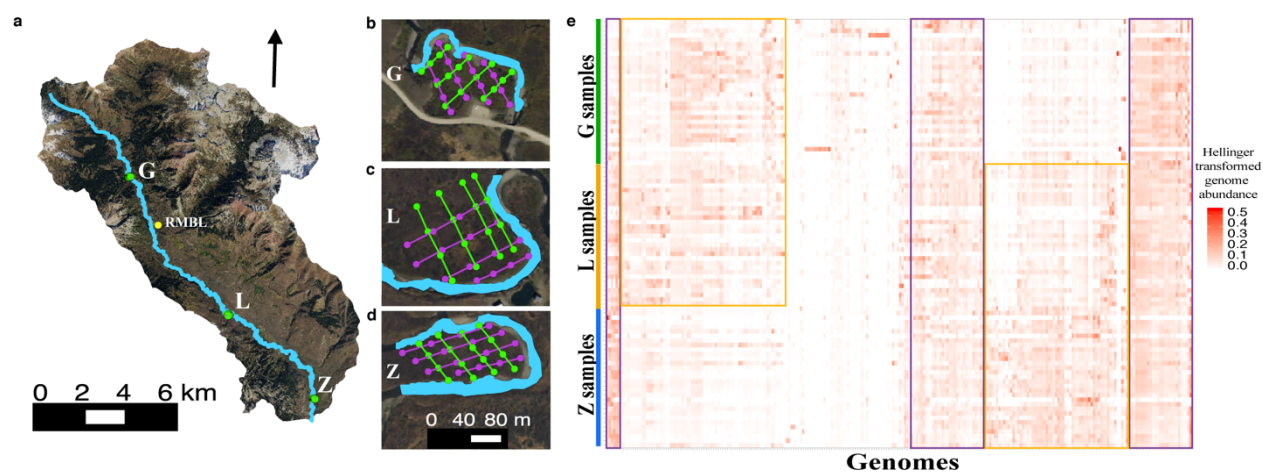
120 We constructed 1,704 draft genomes from three floodplain datasets. About one third (622)
121 of these genomes were classified as draft quality (237 from floodplain G, 150 from floodplain L,
122 and 235 from floodplain Z). After dereplication at 99% average nucleotide identity (ANI) within
123 each floodplain and correction of local assembly errors, we recovered 375 distinct genomes (173
124 from G, 94 from L and 108 from Z). Dereplication across floodplains at 98% ANI generated a
125 final set of 248 representative genomes for further analyses, predominantly $\geq 70\%$ complete
126 ([Supplementary Table 2](#)) and 46% of which were near-complete ($\geq 90\%$).

127 We assessed our genome recovery effectiveness by comparing the number of genomes
128 recovered to a secondary metric for quantifying unique species, the number of unique ribosomal
129 protein L6 (rpL6) marker sequences within unbinned assemblies (see **Methods**). The marker rpL6
130 has been shown to have high recoverability and species delineation accuracy ⁶, relative to methods
131 such as full genome ANI. From the 94 metagenomes, we detected 930 distinct organisms based
132 on rpL6 sequences clustered at 97.5% nucleotide identity. However, 571 of the distinct rpL6
133 sequences were on fragments with coverage that is too low for comprehensive genome sampling
134 ($<7 \times$ coverage given our sequencing depth). The disparity relative to 248 reconstructed
135 representative genomes relative to 359 rpL6 on contigs at $>7 \times$ coverage is attributed to significant
136 challenges associated with genome recovery from soil.

137 Candidate draft genomes were generated for almost all the organisms present at $> 5\text{-}10 \times$
138 coverage in each sample. However, on average, only 5.5% of the total read dataset was stringently
139 mapped (2 mismatches per read of the pair) to the 248 genomes. This is not surprising, given that
140 most sequencing allocations per sample were sufficient to genomically sample only organisms at
141 $> \sim 0.25\%$ relative abundance, and the most abundant organisms in each sample comprised only a
142 few percent of the community.

143 In 2016, we returned to one of the floodplain sites (floodplain L) to collect samples for
144 metatranscriptomics. We performed additional genomic sequencing from 19 of the original 32
145 sites (see **Methods**; [Supplementary Table 1](#)) to provide a reference database for transcript
146 mapping. These new DNA samples were sequenced at an average 3.7 Gbp per sample (2.7 – 4.7
147 Gbp), for a total ~ 0.2 Tbp of sequencing. The RNA samples were sequenced at an average 10.8
148 Gbp per sample (3.2 - 15.1 Gbp) for a total of ~ 0.15 Tbp of sequencing. A total of 299 draft
149 genomes were recovered from these samples, 123 of which passed our quality thresholds after
150 curation. To examine stability across time we pooled the 2015 and 2016 genome sets and
151 dereplicated (at 95% ANI) the combined set of 371 genomes, generating 215 genomes
152 representative of distinct species. Notably, 32 species-level groups were detected in both years and
153 29 were only detected in 2016 ([Supplementary Table 3](#)).

154
155



156
157

158 **Figure 1.** (a) Overview of the East River, CO study site, highlighting the three sampled floodplains (green
159 dots) and the Rocky Mountain Biological Laboratory (RMBL, yellow dot), (b) meander-bound floodplain
160 G, (c) meander-bound floodplain L, (d) meander-bound floodplain Z. Sampling sites as green and purple
161 dots along two sets of four transects. One set of transects in one direction (in green), and the second set of
162 transects along another direction (in purple). (e) Hellinger transformed abundance of dereplicated genomes
163 across samples based on cross-mapping. Genomes and samples clustered by average linkage and Euclidean
164 distance respectively.

165

166

167 *Distribution of organisms within and across meander-bound floodplains*

168 To assess the presence of a representative genome in a sample we relied on the sensitivity
169 of read mapping to the dereplicated genome set. Based on our threshold for detection, about one-

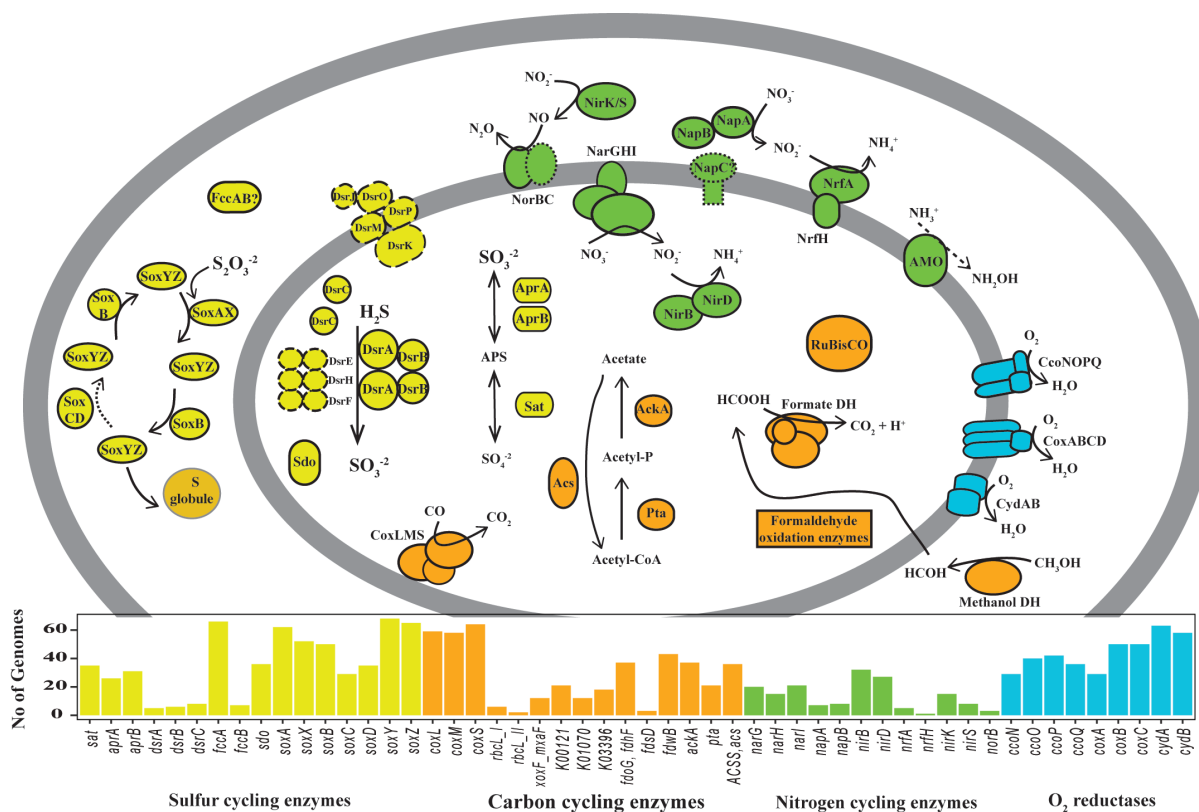
170 third of the genomes were from organisms that were consistently found across floodplains at
171 similar levels of abundance (Purple boxes in **Fig. 1e**). Regardless of their level of abundance, or
172 which floodplain a representative genome was reconstructed from, the genomes were present in a
173 median $> 75\%$ of the samples (78% of upstream floodplain G samples, 84% of mid-stream
174 floodplain L samples, and 87% of downstream floodplain Z samples; [Supplementary Figure 1a](#)).
175 Additionally, the 248 organisms were present in the majority (median 88-91%) of the other
176 samples from the same floodplain from which the genome was reconstructed from ([Supplementary](#)
177 [Figure 1b](#)).

178 Except for some genomes reconstructed from two floodplain G samples, the rest of the
179 genomes were from organisms that shared more similar abundance levels if the floodplains were
180 closer together within the river corridor (Yellow boxes in **Fig. 1e**). More specifically, floodplains
181 G and L or floodplains L and Z shared more organisms than floodplains G and Z, which are located
182 in the upper and lower reaches respectively. Additionally, floodplain G is narrow, and is at times
183 completely flooded, floodplain L is wider and may only flood partially, whereas floodplain Z is
184 the widest and least prone to flooding.

185 Finally, we examined the number of samples where members of a 98% ANI genome cluster
186 were reconstructed from. The 248 genome clusters contained genomes reconstructed from between
187 2 to 39 samples. The largest genome set was for a large group of Betaproteobacteria strains
188 generally related to strains detected in other environments such as soil, sediment and water
189 ([Supplementary Figure 2](#); [Supplementary Data 1](#)). Genomes were reconstructed from two thirds
190 of all samples from floodplain L. This result indicates that strains belonging to this
191 Betaproteobacterial clade may play important roles in floodplain biogeochemistry (**Fig. 2**),
192 especially in soils associated with floodplain L.

193

194



195
196

197 **Figure 2.** Diagram depicting Betaproteobacteria genomes and environmentally relevant capacities encoded
198 by representatives of 98% ANI clusters. Note that no single genome harbors all of these genes, but
199 combinations of them instead (Supplementary Table 4). Some genomes harbor methanol dehydrogenases
200 that are potentially able to turn methanol directly into formate (XoxF type). Enzymes delineated with solid
201 lines were predicted using KOFAM HMMs, and the number of genomes (> 1*) encoding those genes are
202 shown in the bars plot. Enzymes that were predicted using methods as part of ggKbase are shown with
203 dashed lines (long dashes) and enzymes or subunits that are presumably encoded are shown with dotted
204 lines. For more information about metabolic potential see **Methods**. *AMO was included in this diagram
205 even though it was detected in only 1 genome, to indicate aerobic ammonia oxidation is also possibly
206 carried out by members of this clade.

207

208 *Taxonomic composition of the community*

209 Based on the 248 representative genomes detected in each sample, the phylum- or class-
210 level community composition was broadly consistent both within and across floodplains (

216 sample (min = 35 and max = 212). We detected particularly low numbers of genomes in five
217 samples (T157 and T800 from floodplain G and T133, T266 and T620 from floodplain Z),
218 although only samples T133 and T266 from floodplain Z may have been affected by lower
219 sequencing depths ([Supplementary Table 1](#)).

220 Betaproteobacteria was the group with the highest number of representative genomes (80)
221 in all three floodplains. Other abundant taxa across floodplains included Deltaproteobacteria (27
222 representatives), Acidobacteria (21 representatives), Nitrospirae and Planctomycetes (both with
223 13 representatives), Gemmatimonadetes, Gammaproteobacteria, Chloroflexi and Ignavibacteria
224 (12, 11, 11, and 10 respectively).

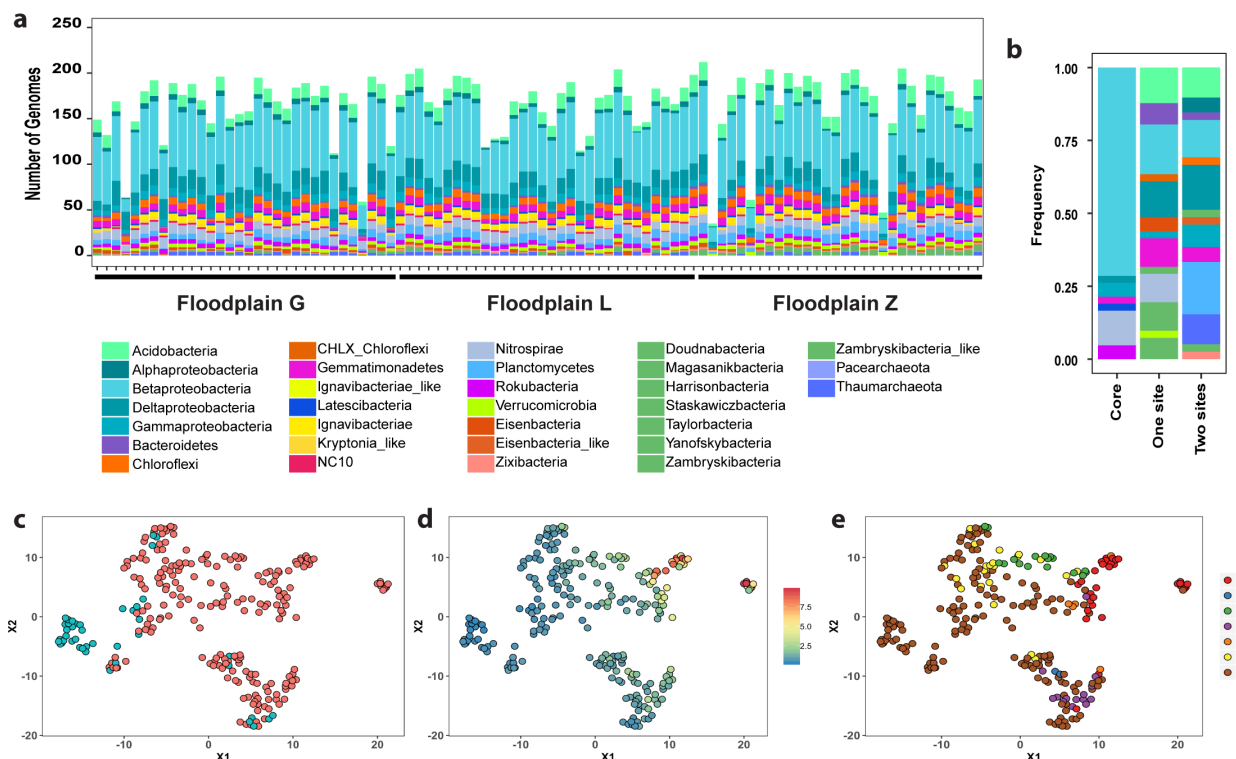
225

226 *Uncovering a core floodplain microbiome*

227 To define a set of organisms representing a core floodplain microbiome we identified
228 organisms that were detected in most sampled sites (≥ 89 of the 94 samples; 90th percentile), and
229 whose abundance did not indicate a statistically significant enrichment in any specific floodplain
230 (by Indicator Species Analysis (ISA); see **Methods** and [Supplementary Table 5](#)). This operational
231 definition resulted in the identification of 42 high prevalence organisms with low variance
232 abundance profiles across all 3 meander-bound sites, which we refer to as the core floodplain
233 microbiome (**Fig. 3c**). In general, genomes with a low coefficient of variation of their abundance
234 (blue dots in **Fig. 3d**) overlapped with genomes that did not display a statistically significant
235 association with any given floodplain (group 7, brown dots in **Fig. 3e**), suggesting a wide
236 distribution of these organisms across floodplains at similar abundance levels. The core floodplain
237 microbiome was dominated by Betaproteobacteria, with lower abundances of Nitrospirae,
238 Rokubacteria, Gemmatimonadetes, Gammaproteobacteria, Deltaproteobacteria, and Candidatus
239 Letescibacteria (**Fig. 3b**).

240 Genomes from organisms not considered to be part of the core floodplain microbiome were
241 associated with one floodplain (G, L, or Z; $n = 41$) or two floodplains ($n = 39$; **Fig. 3e**). Other
242 genomes were not classified as part of the core microbiome because although they were not
243 statistically associated with one or two floodplains, they were not detected in ≥ 89 samples ($n =$
244 126). Genomes affiliated with Acidobacteria, Bacteroidetes, and Chloroflexi were not part of the
245 core floodplain microbiome and were associated with one or two floodplains. The ISA analysis
246 supports the association of some CPR with one floodplain (*i.e.*, between floodplain Z and

247 Yanofskybacteria, Taylorbacteria, Harrisonbacteria, Staskawiczbacteria, and Zambryskibacteria-
 248 like bacteria). Similarly, bacteria in the Verrucomicrobia were associated with one floodplain (Z).
 249 Alphaproteobacteria, Thaumarcheota, Planctomycetes, other CPR (e.g., Zambryskibacteria and
 250 Doudnabacteria) and Eisenbacteria-like bacteria were associated with two floodplains.
 251



252
 253
 254 **Figure 3.** (a) Taxa at the phylum or class level detected across samples and floodplains, samples are in
 255 numerical order (Supplementary Table 1) from the upstream to the downstream floodplain. (b) Taxonomic
 256 composition of genomes in the core floodplain microbiome (core), genomes associated with 1 floodplain
 257 (one site), and genomes associated with two floodplains (two sites). UMAP showing clustering of Hellinger
 258 transformed genome abundances of (c) genomes in the core floodplain microbiome (n = 42; in teal) and
 259 genomes not in the core floodplain microbiome (n = 242; red), (d) and overlay of the coefficient of variation
 260 (ratio of standard deviation to the mean) of genome abundances across samples, and (e) overlay of genomes
 261 associated with individual, pairs, or all floodplains based on an Indicator Species Analysis (ISA). Genomes
 262 that were present in 89 samples or more (teal) were not associated with any particular floodplain by ISA
 263 (group 7 in brown) and their abundance displayed a low coefficient of variation across samples. ISA
 264 genome associations: with floodplain G (1), with floodplain L (2), with floodplain Z (3), with both
 265 floodplain G and floodplain L (4), with both floodplain G and floodplain Z (5), with both floodplain L and
 266 floodplain Z (6), not associated with any particular floodplain (7).

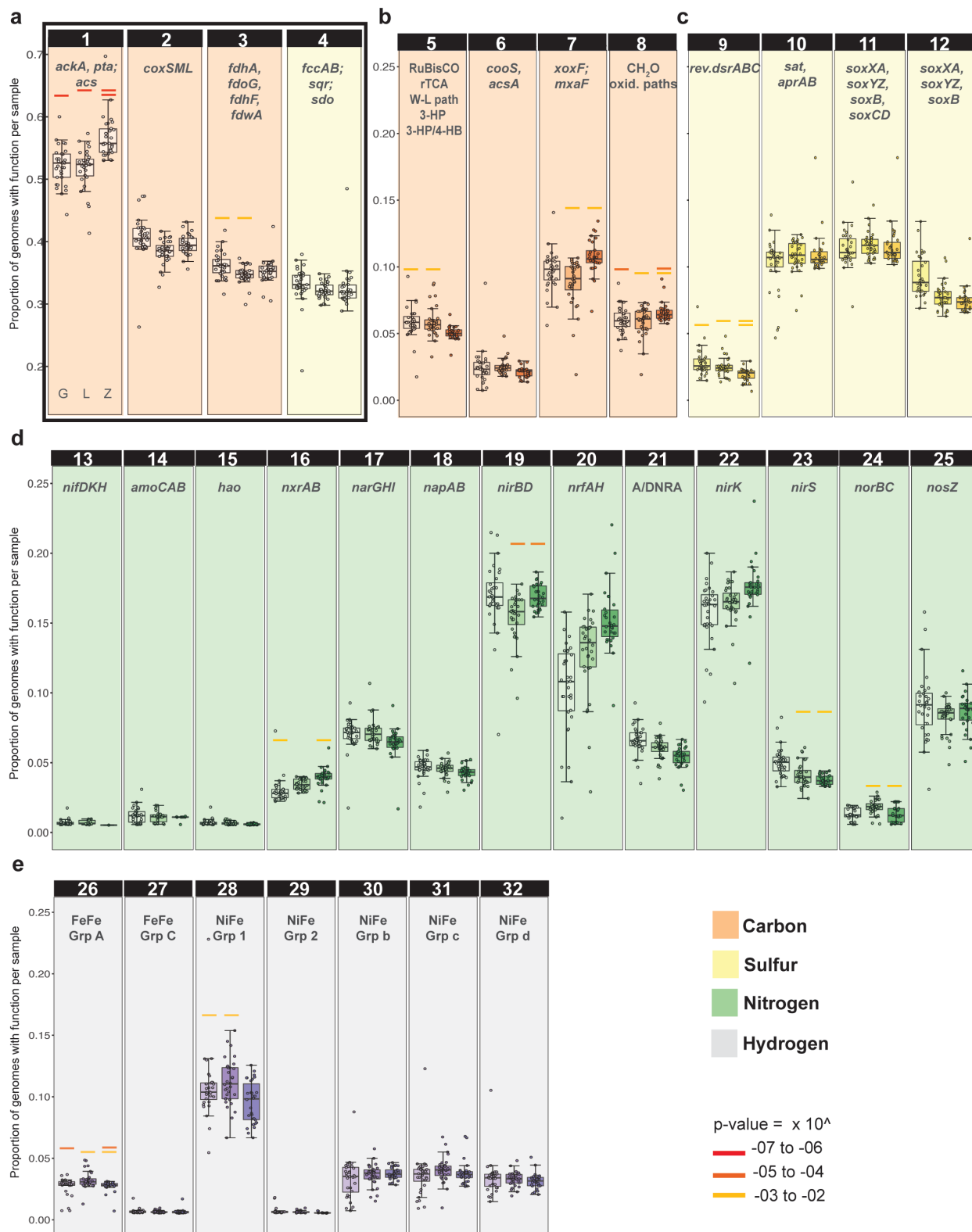
267

268 *Geochemical functions, including those enriched in the core floodplain microbiome*

269 To determine what role floodplain soil Bacteria and Archaea may play in nutrient exports
270 to the East River, we investigated a set of pathways involved in biogeochemical cycling and the
271 microorganisms potentially responsible for them. The biogeochemical processes investigated
272 include oxidation/reduction reactions associated with nitrogen, sulfur and hydrogen, C1 compound
273 metabolism (*e.g.*, CO₂-fixation, CO oxidation, methanogenesis, methane oxidation, methanol
274 oxidation, formate oxidation, methylamine oxidation, formaldehyde oxidation), H₂ consumption
275 or production, and the ability to use O₂ as a terminal electron acceptor for aerobic respiration.

276 A set of HMMs was used to annotate genes encoding for individual protein subunits that
277 make up key enzymes and complete or partial metabolic pathways. For a given ‘function’ (defined
278 as the capacity to carry out a given biogeochemical transformation) to be encoded in a genome,
279 certain criteria for presence had to be met (see **Methods**; [Supplementary Table 4](#)). A total of 32
280 functions comprised the final set of biogeochemical transformations under investigation
281 ([Supplementary Table 6](#)). It is important to note that in some cases we also examined individual
282 steps that are involved in a function, recognizing that some functions could be absent in a single
283 genome because the pathway is carried out by multiple taxa (*i.e.*, steps are encoded in multiple
284 genomes). For example, denitrification occurs in separate steps involving different enzymes, and
285 these steps can be performed by multiple different organisms. Complete ammonia oxidation,
286 anaerobic ammonia oxidation, and methanogenesis (of any kind), were not detected in the
287 dereplicated genome set, although some intermediary steps may still be ecologically relevant.
288 Therefore, some steps involved in these pathways were included in downstream analyses.

289 To study the distribution of the functions of interest among genomes and across
290 floodplains, we determined whether a function was present or absent in each genome in addition
291 to where genomes were detected within and across floodplain samples. To describe the distribution
292 of functions, we calculated the proportion of genomes with a given function compared to the total
293 number of genomes detected in a sample. We found that the ability to use oxygen as an electron
294 acceptor (aerobic respiration) was the most prevalent function among genomes (a median of 70 -
295 85% of genomes in each sample), followed by acetate metabolism (a median of 40 - 65% of
296 genomes in each sample), aerobic carbon monoxide (or other small molecule) oxidation, formate
297 oxidation, and sulfide oxidation (a median of 30 - 50% of genomes in each sample; **Fig. 4a**). This
298 set of functions was consistently present across all three floodplains, whether encoded by the same
299 or different taxa.



300
301
302
303

Figure 4. Proportion of representative genomes at the sub-species level with a function among genomes detected in each sample within each floodplain. In each panel the box plot on the left shows floodplain G,

304 the box plot in the middle shows floodplain L, and the boxplot on the right shows floodplain Z. **(a)** Most
305 abundant functions: **1.** Acetate formation, **2.** Oxidation of CO and other small molecules, **3.** Formate
306 oxidation: CH_2O_2 to $\text{CO}_2 + \text{H}_2$, **4.** Sulfide oxidation: H_2S to S^0 . **(b)** Geochemical transformations in the
307 Carbon cycle: **5.** CO_2 fixation pathways, **6.** Anaerobic CO oxidation, **7.** Methanol oxidation, **8.**
308 Formaldehyde oxidation pathways (see [Supplementary Table 4](#)). **(c)** Geochemical transformations in the
309 sulfur cycle: **9.** Sulfide oxidation (reverse *dsr*) from hydrogen sulfide: H_2S to SO_3^{2-} , **10.** Sulfite oxidation
310 to sulfate (or vice versa): SO_3^{2-} to SO_4^{2-} , **11.** Thiosulfate oxidation without sulfur deposition: $\text{S}_2\text{O}_3^{2-}$ to SO_4^{2-} ,
311 **12.** Thiosulfate oxidation with sulfur deposition: $\text{S}_2\text{O}_3^{2-}$ to $\text{SO}_4^{2-} + \text{S}^0$. **(d)** The nitrogen cycle: **13.** Nitrogen
312 fixation: N_2 to NH_3 , **14.** Ammonia oxidation: NH_3 to NH_2OH , **15.** Hydroxylamine oxidation (requires
313 additional, undetermined enzyme): NH_2OH to NO_2^- , **16.** Nitrite oxidation: NO_2^- to NO_3^- (reversible), **17.**
314 Nitrate reduction (cytoplasmic): NO_3^- to NO_2^- , **18.** Nitrate reduction (periplasmic): NO_3^- to NO_2^- , **19.**
315 Assimilatory nitrite reduction: NO_2^- to NH_4 , **20.** Dissimilatory nitrite reduction: NO_2^- to NH_4 , **21.**
316 Assimilatory or dissimilatory nitrate reduction (ANRA or DNRA): 17 or 18 + 19 or 20, **22 & 23.** Nitrite
317 reduction (Denitrification): NO_2^- to NO , **24.** Nitric oxide reduction: NO to N_2O , **25.** Nitrous oxide reduction:
318 N_2O to N_2 . **(e)** Hydrogen metabolism via hydrogenases: **26.** FeFe hydrogenases group A (fermenting and
319 bifurcating), **27.** FeFe hydrogenases group C (H_2 sensors), **28.** NiFe hydrogenases group 1 (H_2 oxidation),
320 **29.** NiFe hydrogenases group 2 (H_2 oxidation), **30.** NiFe hydrogenases group 3b (bidirectional), **31.** NiFe
321 hydrogenases group 3c (bidirectional), **32.** NiFe hydrogenases group 3d (bidirectional). Paired colored bars
322 above any two given boxplots with the same color and at the same level indicate statistically significant
323 differences between those two floodplains (two-way ANOVA).
324

325 We also considered the distribution of functions that were detected in < 25% of genomes
326 (**Fig. 4b-e**). Of the remaining C1 transformations examined, methanol oxidation to formaldehyde
327 was found in a median of ~10% of the genomes. Of the sulfur transformations, sulfite (SO_3^{2-})
328 oxidation to sulfate (SO_4^{2-}), and thiosulfate ($\text{S}_2\text{O}_3^{2-}$) oxidation without sulfur (S^0) deposition and
329 thiosulfate oxidation with sulfur deposition were most prevalent. For reactions involving hydrogen
330 consumption or formation, genes encoding group 1 NiFe hydrogenases (likely used for H_2
331 oxidation) were found in a higher proportion of genomes than any other types of hydrogenases.

332 Nitrogen transformations were studied individually and as part of the nitrogen cycle. We
333 found the capacity for nitrate (NO_3^-) use as a terminal electron acceptor in dissimilatory NO_3^-
334 reduction in a substantially lower proportion of genomes (2 - 10%) than the capacity to use O_2 as
335 a terminal electron acceptor. Of the reactions involved in nitrification, namely ammonia oxidation,
336 hydroxylamine oxidation and nitrite (NO_2^-) oxidation, genomes encoding the oxidation of nitrite
337 via nitrite oxidoreductase (NXR) were more common than genomes encoding the first two steps.
338 The capacity for NO_3^- reduction as part of denitrification (NapAB or NarGHK) was encoded by
339 far fewer genomes than NO_2^- reduction (which can be carried by via multiple enzymes, including
340 NirK, NirS, NrfAH for dissimilatory nitrite reduction or NirBD for assimilation). Fewer genomes

341 are predicted to encode the capacity to reduce nitric oxide (NO, the product of nitrite reduction) to
342 nitrous oxide (via NorBC) than genomes with the capacity for nitrous oxide (N₂O) reduction to
343 N₂. Overall, the most prevalent genomically encoded function was nitrite reduction, and capacities
344 for consecutive nitrogen cycling steps were typically encoded in multiple different genomes. In
345 other words, there is evidence to support the prevalence of metabolic handoffs ⁷ in the nitrogen
346 cycle.

347 We identified functions that were significantly enriched (FDR ≤ 0.05; hypergeometric test)
348 in the core floodplain microbiome (a subset of ISA group 7) and found that the capacities to use
349 O₂ as a terminal electron acceptor, to perform aerobic CO or other small molecule oxidation, and
350 thiosulfate oxidation (both with and without sulfur deposition) were enriched in these organisms.

351

352 *Environmental factors as drivers of function distribution across and within floodplains*

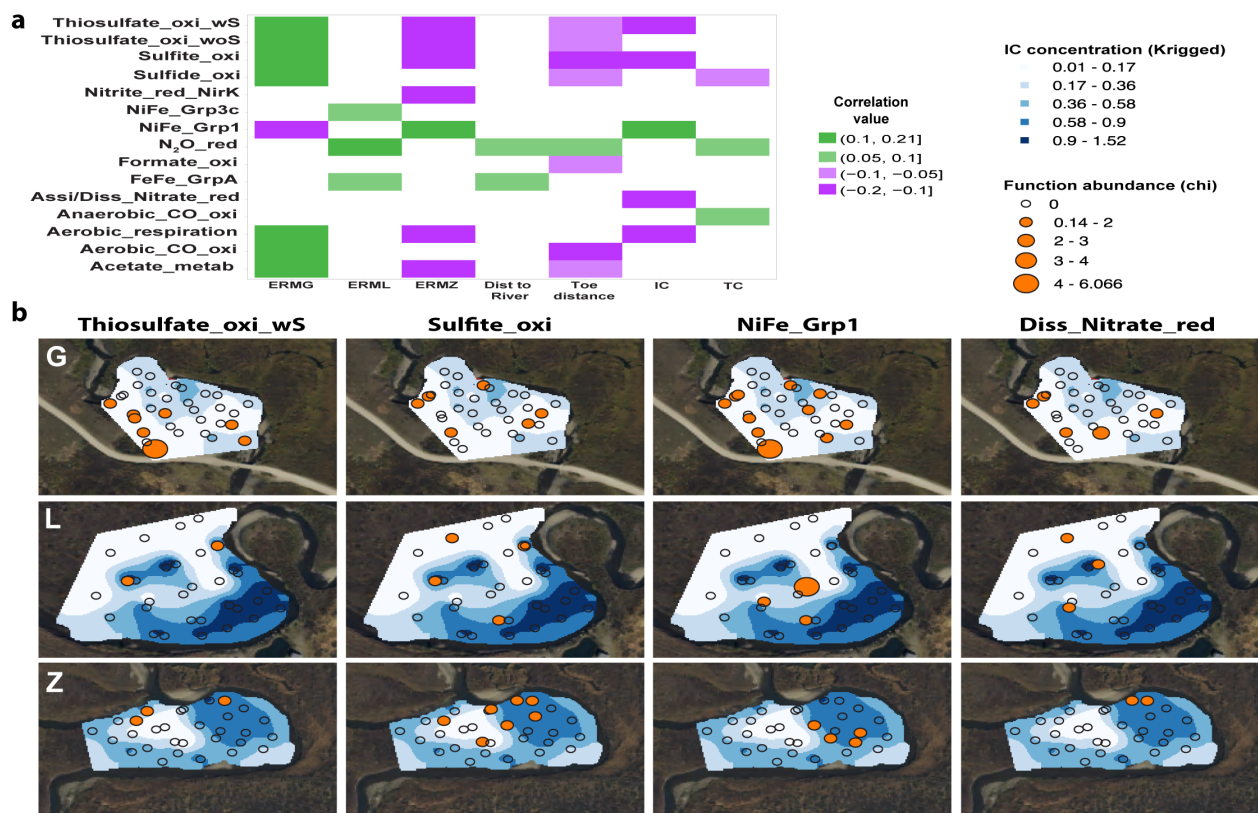
353 Environmental variables ([Supplementary Figure 4](#)) may explain in part the patterns of
354 enrichment of genomically encoded functions described above. We first looked into correlations
355 involving the following variables: total carbon (TC), total organic carbon (OC), total inorganic
356 carbon (IC), total nitrogen (TN), organic carbon to nitrogen ratio (OC:N), distance of a sample to
357 the river (Dist. to river), easting and northing (cartesian coordinates for position on the floodplain),
358 distance to the inner bank edge (from here on: toe distance) and distance to middle of the meander-
359 bound floodplain as alternative measures of position on the floodplain ([Supplementary Figure 5](#)),
360 topographic position index (TPI; as a proxy for the likelihood a site would be flooded during
361 periods of high discharge or snowmelt), and elevation. Statistical analysis indicated that TC, OC,
362 TN, and OC:N were all highly correlated with each other across the same set of metagenomic
363 samples ([Supplementary Figure 6](#)) and their individual effects were not possible to disentangle.
364 Thus, we chose either TC or OC for downstream analyses. Given the Northwest to Southeast
365 orientation of the watershed, elevation, Easting, and Northing were all highly correlated with
366 floodplain (G vs L vs Z), so only floodplain was included as a categorical variable. In summary,
367 TC, floodplain, IC, TPI, distance to the river and toe distance were the variables evaluated with
368 the fourth corner method ⁸ to assess the response of each function at the gene level to the selected
369 environmental or soil chemistry and GIS variables (see **Methods**).

370 A group of biogeochemical transformations (gene level) displayed some correlation with
371 environmental variables, particularly with individual floodplains (**Fig. 5a**). Genome abundances

372 were used as proxies for abundance of functions each genome encoded. The upstream floodplain
 373 G was positively correlated with thiosulfate oxidation (with and without S deposition), sulfite
 374 oxidation, sulfide oxidation, O₂ as a terminal electron acceptor, aerobic CO or other small molecule
 375 oxidation, and acetate metabolism. Only N₂O reduction was positively correlated with the middle
 376 floodplain L. The downstream floodplain Z was positively correlated with H₂ oxidation via group
 377 1 NiFe hydrogenases (a function that was negatively correlated with upstream floodplain G). Most
 378 sulfur compound transformations, as well as nitrite reduction, aerobic respiration and acetate
 379 metabolism were negatively correlated with floodplain Z (**Fig. 5a**).

380 Overall, genomes with the capacity for aerobic respiration and sulfur compound oxidation
 381 are most prevalent towards the headwaters (floodplain G), and within this floodplain sulfur
 382 compound oxidation apparently is associated with low IC (**Fig. 5b**). Within the downstream
 383 meander where aerobic respiration is least prominent in the genomes, bacteria able to oxidize H₂
 384 via Group 1 NiFe hydrogenases appear correlated with somewhat elevated concentrations of IC
 385 (**Fig. 5b**).

386



387
 388

389 **Figure 5.** Function abundance and its correlation with environmental variables. **(a)** Significant positive
390 (green) or negative (violet) correlations between environmental variables (bottom) and biogeochemical
391 transformations (left) identified by a fourth corner analysis. **(b)** Abundance of genomes encoding functions
392 positively correlated with inorganic carbon concentrations (IC; %).

393

394 *Potentially active genes encoding key biogeochemical transformations in the riparian zone*

395 To determine whether key functions encoded in the genomes were transcriptionally active
396 at the time of sampling (early September 2016; during base flow conditions like previous year),
397 we re-sampled floodplain L for metatranscriptomics and metagenomics. This floodplain was
398 chosen among the three because it shared the majority of organisms detected in 2015 with the other
399 two floodplains.

400 Considering potential differences between the two years, metatranscriptomic reads were
401 mapped to a dereplicated genome set at the species level (95% ANI), which comprised 215
402 genomes reconstructed from samples collected in 2015 and 2016. We calculated transcript counts
403 using read pairs mapped to predicted open reading frames (ORFs) with at least 95% nucleotide
404 identity (see **Methods**). The highest median transcript counts were observed for Nitrospirae and
405 Betaproteobacteria, followed by Candidatus Latescibacteria and Eisenbacteria-like bacteria,
406 Rokubacteria, and Deltaproteobacteria (**Fig. 6a**). We also evaluated the number of reads mapping
407 to genes encoding key functions and determined what percentile in the distribution of transcription
408 levels each gene fell in.

409 Key genes involved in potentially active biogeochemical transformations with a median
410 transcription > 75th percentile of all the genes transcribed in a given genome included *amoCAB*,
411 and *nxrAB*, involved in nitrification. The *amoCAB* genes for aerobic ammonia monooxygenase
412 were found in the 90th percentile of transcribed genes across genomes. However, these genes were
413 present in very few genomes (one Nitrospirae and two Thaumarcheota). Similarly present in few
414 genomes and also highly transcribed were genes involved in CO₂ fixation, specifically RuBisCO
415 forms I and II and enzymes in the reductive TCA cycle (OFOR and citrate lyase). Other highly
416 transcribed genes were for methanol oxidation to formaldehyde (*xoxF*, *mxoF*) and formate
417 oxidation (*fdhAB*, *fdoG*, *fdhF*, *fdwA*, *fdsD*, *fdwB*) as part of C1 metabolism.

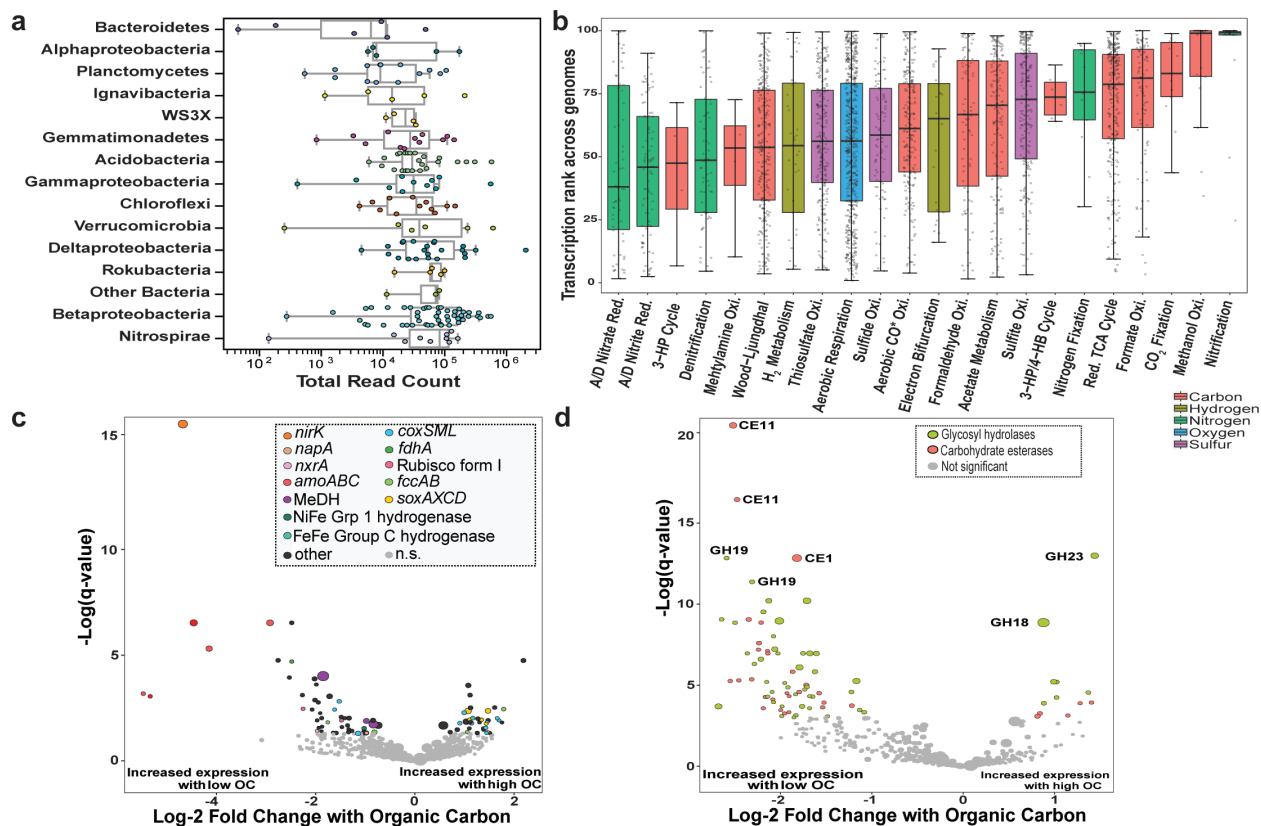
418 Functions enriched in the core floodplain microbiome (aerobic CO or other small molecule
419 oxidation, thiosulfate oxidation, and the ability to use O₂ as a terminal electron acceptor via
420 *coxABCD*, *cydAB* or *ccoN*) and functions that displayed some degree of correlation with
421 environmental variables in gene abundance (e.g., sulfite oxidation via *sat* and *aprAB* or *dsrAB* and

422 sulfide oxidation via *fccAB*) were most often between 50th-75th percentile of transcribed genes per
423 genome. Surprisingly given their prominence in genomes, genes involved in nitrogen cycling such
424 as *narGHI* or *napA* and especially *nrfAH* responsible for dissimilatory nitrite reduction, *nirK* and
425 *nosZ* responsible for some denitrification steps, displayed transcription levels only between the
426 35th - 50th percentiles (**Fig. 6b**).

427 We tested for differential transcription levels in response to changes in environmental
428 variables ([Supplementary Table 7](#)) using DESeq2⁹. Of the four environmental variables that were
429 highly correlated with each other (TC, OC, TN, OC:N; either positively or negatively
430 [Supplementary Figure 7](#)), we observed the strongest differential gene expression in response to
431 OC (**Fig. 6c**). In samples with higher concentrations of OC, genes involved in the Sox pathway for
432 thiosulfate oxidation (*soxAX* and *soxCD*) and those involved in aerobic CO or other small molecule
433 oxidation (*coxLMS*) were highly transcribed. More specifically, transcripts mapped to one *coxL*
434 form I gene (true carbon monoxide dehydrogenase, CODH) and the rest mapped to four other *coxL*
435 form II genes (carbon monoxide-like dehydrogenase). The form I transcripts were correlated with
436 high OC, and the form II transcripts with both low (1 hit) and high OC (3 hits).

437 In samples with low OC, highly transcribed genes included those involved in methanol
438 oxidation (*xoxF*, *mxoF*), formate oxidation (*fdhA*), sulfide oxidation (*fccAB*), hydrogen
439 metabolism (NiFe Grp1), ammonia oxidation (*amoABC*), nitrite oxidoreduction (*nxrA*), nitrate
440 reduction (*napA*), and nitrite reduction (*nirK*). Samples with low OC also have low TN, which
441 may in part be attributed to high activity of microbial nitrification followed by denitrification, with
442 denitrification reliant on consumption of OC.

443 As might be expected, RuBisCO form I was highly transcribed under conditions of low
444 OC. Activity of the Calvin Benson Bassham (CBB) pathway for CO₂ fixation is linked to
445 Betaproteobacteria, Deltaproteobacteria, Gammaproteobacteria, and NC10, some of which have
446 metabolisms fueled by oxidation of intermediate sulfur compounds (e.g., sulfide or thiosulfate
447 oxidation). The observations reveal a potentially important source of organic carbon in some
448 floodplain soils. In terms of overall transcriptional activity, autotrophic pathways may not be
449 expressed but the organisms may otherwise be highly transcriptionally active. In fact, mostly
450 organisms from the phyla Thaumarchaeota, Rokubacteria, NC10 and Nitrospirae were active under
451 low OC conditions. Betaproteobacteria and Acidobacteria were transcriptionally active under
452 conditions of higher OC.



453
 454

455 **Figure 6.** Analyses of transcription from samples collected in 2016 mapped to the species level
 456 representative genome set. **(a)** Transcription activity of genomes grouped by phylum or class based on total
 457 transcript read counts mapped to the representative genomes. Other Bacteria: Candidatus Latescibacteria
 458 and Eisenbacteria-like bacteria. Phylogeny of the genomes at the species level was confirmed based on a
 459 concatenated ribosomal proteins tree (Supplementary Data 2). **(b)** Average transcription percentile for all
 460 genes encoding enzymes involved in a given biogeochemical transformation in each representative genome
 461 across all 2016 metatranscriptomes. **(c)** Differentially transcribed genes in response to soil OC. Statistically
 462 significant (DESeq2; $q < 0.05$) genes are colored by function and not significant (n.s.) genes are in grey.
 463 **(d)** Differentially transcribed genes encoding CAZY in response to soil OC.

464

465 We also investigated the potential for organic matter degradation through transcription of
 466 genes encoding carbohydrate-active (CAZY) enzymes (Supplementary Table 8)¹⁰. We narrowed
 467 our search to CAZY enzyme types that were present in at least 60% of the genomes, as a proxy for
 468 widespread distribution in the floodplain soil microbial community. The most abundantly
 469 transcribed CAZY genes were in the glycosyl hydrolase (GH) and carbohydrate esterase (CE)
 470 classes. In general, the highly transcribed enzymes in the CE class use hemicellulose and amino
 471 sugars as substrates, resulting in acetate as a byproduct. Acetate could be utilized by many
 472 floodplain organisms, considering the prevalence of genes involved in acetate metabolism among

473 the genomes. Enzymes in the GH class use cellulose, pectin, chitin and starch as substrates,
474 releasing a variety of sugars as byproducts, which can be utilized for central metabolism during
475 growth. We then tested for CAZY differential expression in response to changing concentrations
476 of organic carbon. 12 CAZYS (both GH and CEs) expressed by three strains of Betaproteobacteria
477 increased in expression in samples with high OC. Many of the same classes of GH and CE
478 displayed high levels of transcription correlated with both low and high OC levels (*e.g.*, GH23,
479 GH28, CE4, and CE11 **Fig. 6d**), although gene expression by three strains of Betaproteobacteria
480 correlated with high levels of OC. Similar CAZY enzymes were expressed by NC10, Nitrospirae,
481 and Rokubacteria, the same organisms commonly associated with high levels of transcription
482 under low OC.

483

484 **Discussion**

485 Biogeochemical processes modulate C, S, and N exports from watersheds, including the
486 East River ¹¹. Important questions relate to the sources and sinks of these compounds and the
487 biological controls on them. Some data indicate that a subset of the organic carbon in sediments
488 from East River floodplains derives from the shale, although plants are the obvious central source
489 for fixed carbon in areas of more developed soils ¹². CO₂ fixation genes were relatively rarely
490 detected in the bacterial genomes, which might be interpreted to support this deduction. However,
491 genes for CO₂ fixation in a few organisms were very highly transcribed, indicating at least periodic
492 inputs of microbially-produced organic carbon into riparian zone soils. Spatially, high activity of
493 genes involved in CO₂ fixation was correlated with low organic carbon concentrations in soil.
494 Many organisms predicted to rely on CO₂ fixation as their main carbon source are aerobic
495 chemolithoautotrophs that oxidize inorganic compounds (*e.g.*, NH₃⁺, NO₂⁻, S⁰, H₂S, H₂, CO, S₂O₃)
496 as a source of energy. Thus, we infer significant linkages amongst these key element nutrient
497 cycles.

498 Low concentration of organic carbon also correlated with high activity of genes involved
499 in methanol oxidation. Methanol results from the breakdown of plant material, such as pectin and
500 lignin, and the activity of methanol dehydrogenases may be indicative of decomposed organic
501 matter. Similarly, low concentration of organic carbon correlated with high activity of genes
502 involved in sulfide and H₂ oxidation, nitrification and interconversion of nitrite and nitric oxide.
503 The organisms responsible for these reactions are primarily autotrophs.

504 Interestingly, the capacities for thiosulfate oxidation/elemental sulfur formation, sulfite
505 oxidation and H₂ oxidation, as well as assimilatory or dissimilatory nitrate reduction to ammonia
506 (ANRA or DNRA) were patchily spatially distributed (**Fig. 4b**), possibly localized by lower
507 inorganic carbon concentrations. Further, genes for sulfur compound oxidation were more
508 prominent in genomes of organisms from the upstream floodplain, which is closer to the adjoining
509 hill, possibly reflecting higher inputs of intermediate sulfur compounds from rock weathering
510 reactions in the headwater compared to downstream regions. The sources of thiosulfate could be
511 weathering of detrital grains of shale-associated pyrite and/or reoxidation of microbially-produced
512 sulfide in anoxic OC-rich regions of the soil or underlying river sediment. Closer proximity to
513 igneous intrusives in the upstream part of the drainage (*e.g.*, near Floodplain G) leads to greater
514 incidence of pyrite-bearing shales. It has been shown previously that hydrological connectivity can
515 shape microbial activity, with low connectivity linked to higher abundance of genes involved in
516 sulfur metabolism¹³. In the East River, sulfur compounds may be redistributed from upstream to
517 downstream regions, but the degree of hydrologic connectivity within and across floodplains is
518 uncertain and varies dramatically over the course of the year¹⁴. Additionally, these shallow soils
519 may only be hydrologically connected to the river during high water flood events or through
520 vertical transport.

521 By contrast, high OC levels correlated with high activity of genes involved in oxidation of
522 CO (form I CODH), and other small carbon compounds (form II or other subtypes), which may
523 be substrates for carbon monoxide dehydrogenases⁵. CO may be sourced from the atmosphere, by
524 thermochemical, photochemical, and chemical degradation of organic matter in soils and marine
525 sediments, and from biological production by microbes, leaves, roots and animals¹⁵. East River
526 CO oxidizers are most likely carboxydovores that require organic carbon to grow, even though
527 they can oxidize CO at atmospheric levels (*i.e.*, they use a high affinity form I CODH)^{16,17}. This
528 is in contrast to carboxydotrophs that grow with CO as the sole energy and carbon source and
529 require CO at greater than atmospheric concentrations (for a low affinity form I carbon monoxide
530 dehydrogenase (CODH);¹⁵). Additionally, form II CO dehydrogenases seem to play a key role in
531 this ecosystem, although very little is known about their actual function. Detection of a high
532 prevalence of CODH and CODH-like enzymes echoes results from a grassland soil system
533 Diamond, et al.⁵, reinforcing the suggestion that small carbon compounds such as plant exudates,
534 may be an important carbon currency under some conditions.

535 High organic carbon levels also correlated with high activity of genes involved in
536 thiosulfate oxidation, and carbohydrate esterases and glycosyl hydrolases such as GH23
537 (lysozyme) and GH18 (chitinase). These GHs would be required for organic matter degradation at
538 locations of higher carbon availability, where presumably plants and fungi are more abundant.
539 Bacteria that degrade plant biomass are also known to employ catabolite repression of CAZy
540 enzymes^{18,3517}, perhaps explaining the lower diversity of CAZys under these conditions. Different
541 variants of these carbohydrate-active genes were highly expressed in a variety of taxa including
542 Rokubacteria, Nitrospirae and NC10 in soil with low organic carbon, where diverse carbon sources
543 must be exploited for survival.

544 Many watershed ecosystems are limited by access to biologically available nitrogen, the
545 important sources of which are likely to be shale bedrock weathering¹⁹, atmospheric deposition
546²⁰, and nitrogen fixation. A complex interplay of biological processes impact nitrogen speciation
547 and bioavailability, including ammonia oxidation (nitrification), denitrification to N₂, and nitrite
548 assimilation via ANRA or DNRA. The nitrogen budget can be addressed by direct measurement
549 of inputs, plant-associated inventories, and the concentration of inorganic and organic nitrogen
550 compounds exported from the watershed via rivers^{21, 22}. By comparing these numbers, it may be
551 possible to estimate the fraction of the bioavailable nitrogen that is lost from the system via loss
552 as N₂ and trace gases. What is missing from this analysis is an estimate of the degree to which
553 nitrogen compounds are recycled, the role of riparian zone soils in these processes, and the
554 potential for subsurface storage of nitrogen compounds in microbial biomass.

555 Using genome-resolved metagenomics we identified the capacity for nitrogen fixation and
556 ammonia oxidation to nitrite and nitrate in relatively few organisms, yet the metatranscriptomic
557 data show these to be highly active functions. Thus, we infer important microbial contributions to
558 reservoirs of oxidized nitrogen compounds in riparian zone soils, with the potential to substantially
559 augment inputs from atmospheric deposition and bedrock weathering. Genes involved in nitrite
560 reduction (dissimilatory nitrate reduction to ammonium or denitrification to N₂), while abundant
561 in comparison to other capacities for nitrogen transformation, displayed surprisingly low levels of
562 transcription at the time of sampling. Over the course of the year, the fluctuating water table and
563 periodic flooding should provide environmental niches for both obligately aerobic, and anaerobic
564 processes. Furthermore, the soil oxidation state and the carbon to nitrate ratio, particularly in
565 nitrogen-limited systems, may favor DNRA over denitrification²³. The current study was

566 conducted during base flow conditions (both years), well after the snowmelt season, when high
567 river discharge induces flooding and therefore anoxic conditions. In the meander-bound
568 floodplains, snowmelt-derived flow in this ecosystem persists well into the year¹², so shallow soils
569 may be flooded long after discharge levels drop. The results raise the possibility of coupling of
570 nitrification and dissimilatory nitrate pathways on a temporal basis, under baseflow conditions
571 (when nitrification is dominant), or under snowmelt conditions (when dissimilatory processes
572 occur). High net nitrification has been reported for riparian zones when the water table is below -
573 30 cm²⁴, accordingly the water table in floodplain L was observed to be below this level in
574 September 2016. Overall, re-assimilation of nitrogen as ammonium may be important in this
575 ecosystem, particularly if nitrogen limited.

576 An important result from the current study was that there appears to be a core floodplain
577 microbiome composed of specific bacterial species from Betaproteobacteria,
578 Gammaproteobacteria, Deltaproteobacteria, Nitrospirae, Candidatus Latescibacteria, and
579 Rokubacteria; and all of these groups were transcriptionally active at the time of sampling. Many
580 of the clusters of related genomes are relatively distantly related to previously described bacterial
581 types. Thus, we conclude that many of the most abundant taxa in these riparian zone soils are
582 organisms that have, until now, remained essentially outside of the range of scientific
583 investigations. Importantly, capacities for aerobic respiration, aerobic oxidation of CO and other
584 small molecules, as well as thiosulfate oxidation with formation of elemental sulfur, were enriched
585 in the core floodplain microbiome. Notably, the most abundant functions of the core microbiome
586 were only moderately transcribed at the time of sampling.

587 In general, we found that gene and organism abundances do not predict transcription levels.
588 The *in situ* transcription data revealed the potentially very high importance of rare genes and
589 organisms. However, it is important to note that transcript datasets are a snapshot from a moment
590 in time, and that transcription patterns will vary across seasons and maybe even daily. Notably,
591 our analyses showed organismal and functional overlap in microbial communities found both
592 within and across the three floodplains over two consecutive years (~15% of species in common,
593 despite the very high diversity and complexity of the soils). Thus, in contrast to potentially
594 substantial transcriptome variability, gene inventories reflect metabolic potential that likely
595 remains fairly constant throughout the year. Thus, we conclude that, at the watershed scale,
596 meander-bound regions of floodplain soils are “functional zones” that likely predict

597 biogeochemical transformations along the riparian corridor, thereby providing broadly
598 generalizable inputs to ecosystem models.

599

600 **Methods**

601

602 *Study site and samples collection*

603 The East River (ER) watershed has been described elsewhere³. In brief, the ER watershed
604 is a 300 km² area largely underlain by marine shales of the Cretaceous Mancos formation located
605 in the Elk Mountains in west-central Colorado. The ER is a headwaters catchment in the Upper
606 Colorado River basin, with an average elevation of 3350 m. At about 62 km long, the ER traverses
607 an elevational gradient that includes alpine, subalpine, and montane life zones as a function of
608 stream reach. The average annual temperature is ~ 0 °C, with long cold winters and short cool
609 summers, and the majority of precipitation is received in the form of snow²⁵.

610 The sampling sites are located across an altitudinal gradient followed by the river (~2700
611 – 2900 m). The floodplain at the highest elevation is located ca. 6 km from the headwaters, nearby
612 Gothic, Colorado, site of the Rocky Mountain Biological Laboratory (RMBL, **Fig. 1**). Therefore,
613 samples collected from this site were named East River Meander-bound floodplain G (ERMG).
614 The second site was located ca. 8 km downstream of Gothic, among a series of floodplains, one of
615 which is situated adjacent to an intensive research site of the Watershed Function SFA³. This
616 floodplain stands out because of its larger size, and samples were named ERML (L for large). The
617 third site was located ca. 18 km downstream of Gothic and just upstream of the confluence with
618 Brush Creek. Samples from this site were named ERMZ, with the stream reach between ERML
619 and ERMZ being characterized by a relatively low gradient with high sinuosity.

620 In September 2015, during base flow conditions, two series of perpendicular transects were
621 laid out at each site. Each set of transects comprised four transects that were parallel between them
622 (**Fig. 1**). One set of transects were approximately North to South (T1-T4) and the other set of
623 transects were East to West (T5-T8). The starting point of each transect was designated “0 m” and
624 the location of the other sites along the transect was relative to the point of origin. A Trimble Geo
625 7X GPS was used to determine the exact location of each site along the transects with an accuracy
626 of 0.5 m. The distance (in meters) of each sample to the point of origin was included in the sample
627 name, which comprised the initials of the study area (ER), the initials for each meander-bound

628 floodplain (*i.e.*, MG, ML or MZ), the transect number (*i.e.*, T1-T8), and the distance in meters
629 from the first sample collected at the start point (*e.g.*, 19 m). We sampled an area ~ 4,600 m² in
630 floodplain G, ~ 8,000 m² in floodplain L, and ~ 5,400 m² in floodplain Z.

631 Four soil samples from the 10-25 cm (\pm 1-2 cm) soil depth interval were collected in the
632 span of 10 days along each one of the eight transects, for a total of 32 samples per floodplain. Each
633 site was cleared of grasses and other vegetation with clippers, and the first ~ 10 cm of soil was
634 removed with a sterile shovel. Soil samples were collected using sterile tools, including a soil core
635 sampler and 7.6 x 15.2 cm plastic corer liners (AMS, inc), stainless-steel spatulas, and Whirl-pak
636 bags. Samples were immediately stored in coolers for transportation to RMBL, where samples
637 were prepared for archival and transportation to the University of California, Berkeley. Soil cores
638 were broken apart and manually homogenized inside the Whirl-pak bags. Subsamples for chemical
639 analyses, DNA extractions, and long-term archival were obtained inside a biosafety cabinet, kept
640 at – 80 °C, transported in dry ice, and stored at – 80 °C at the University of California, Berkeley.

641 In September 2016, another round of samples collection was conducted at floodplain L for
642 metagenomics, metatranscriptomics, and chemical analyses. A subset of 19 out of the 32 sampling
643 sites from the previous year was targeted, and a subset (15) of those was also selected for
644 metatranscriptomics ([Supplementary Table 1](#)). Given that floodplain L was the site with the lowest
645 total number of draft genomes recovered in 2015, we added new sites closer to the original sites
646 with the intent of increasing this number by leveraging differential coverage across samples ²⁶.
647 Four new sites located in between the original transects (denominated ERMLIBT) and two sites
648 adjacent to ERMLT660 (ERMLT660_1 and ERMLT660_2) were sampled. Additionally, samples
649 were collected from above the water table (approximately below 40-50 cm from the surface) at a
650 depth of 32-47 cm (\pm 4-6 cm) from three sites (ERMLT200, ERML231 and ERML293) along T2.
651 Samples from the 11-25 cm (\pm 1-1 cm) soil layer were obtained following the same protocol as
652 the previous year, with the exception that subsamples for RNA sequencing were preserved *in situ*.
653 Once the soil cores were transferred to a Whirl-pak bag, they were manually homogenized inside
654 the bags. Eight grams (8 g) of soil were collected using sterile stainless-steel spatulas directly into
655 50 mL sterile falcon tubes containing 20 mL of LifeGuard Soil Preservation Solution (formerly
656 MoBio) for RNA preservation. The samples were mixed by hand to saturation with the LifeGuard
657 solution, stored in a chilled cooler for transportation to RMBL and later stored at -80 °C.
658

659 *Soil chemistry*

660 Total carbon (TC) and total inorganic carbon (TIC) were analyzed using a Shimadzu TOC-
661 VCPH analyzer equipped with a solid sample module SSM-5000A (Shimadzu Corporation,
662 Japan). Total organic carbon (TOC) was obtained from the difference between TC and TIC. For
663 TC quantification, a subsample of the dried solids was weighed into a ceramic boat and combusted
664 in a TC furnace at 900 °C with a stream of oxygen. To ensure complete conversion to CO₂, the
665 generated gases are passed over a mixed catalyst (cobalt/platinum) for catalytic post-combustion.
666 The CO₂ produced is subsequently transferred to the NDIR detector in the main instrument unit
667 (TOC-VCSH). Quantification of the inorganic carbon was carried out in a separate IC furnace of
668 the module. Phosphoric acid is added to the sample and the resulting CO₂ is purged at 200 °C and
669 measured.

670 Total nitrogen (TDN) was analyzed using a Shimadzu Total Nitrogen Module (TNM-1)
671 coupled to the solid sample module (SSM-5000A) and TOC-VCSH analyzer (Shimadzu
672 Corporation, Japan). TNM-1 is a non-specific measurement of TN. All nitrogen species in samples
673 were combusted at 900 °C, converted to nitrogen monoxide and nitrogen dioxide, then reacted
674 with ozone to form an excited state of nitrogen dioxide. Upon returning to the ground state, light
675 energy is emitted. Then, TN is measured using a chemiluminescence detector.

676

677 *DNA extraction and sequencing*

678 Genomic DNA was extracted from ~10 g of thawed soil using Powermax Soil DNA
679 extraction kit (Qiagen) with some minor modifications as follows. Initial cell lysis by vortexing
680 vigorously was substituted by placing the tubes in a water bath at 65 °C for 30 minutes and mixing
681 by inversion every 10 minutes to decrease shearing of the genomic DNA. After adding the high
682 concentration salt solution that allows binding of DNA to the silica membrane column used for
683 removal of chemical contaminants, vacuum was used instead of multiple centrifugation steps.
684 Finally, DNA was eluted from the membrane using 10 mL of the elution buffer (10 mM Tris
685 buffer) instead of 5 mL to ensure full release of the DNA. DNA was precipitated out of solution
686 using 10 mL of a 3 M sodium acetate (pH 5.2) and glycogen (20 mg/mL) solution and 20 mL
687 100% sterile-filtered ethanol. The mix was incubated overnight at 4 °C, centrifuged at 15,000 x g
688 for 30 minutes at room temperature, and the resulting pellet was washed with chilled 10 mL sterile-
689 filtered 70% ethanol, centrifuged at 15,000 x g for 30 min, allowed to air dry in a biosafety cabinet

690 for 15-20 minutes, and resuspended in 100 μ L of the original elution buffer. Genomic DNA yields
691 were between 0.1 – 1.0 μ g/ μ L except for two samples with 0.06 μ g/ μ L. Power Clean Pro DNA
692 clean up kit (Qiagen) was used to purify 10 μ g of DNA following manufacturer's instructions
693 except for any vortexing was substituted by flickering of the tubes to preserve the integrity of the
694 high molecular weight DNA. DNA was resuspended in the elution buffer (10 mM Tris buffer, pH
695 8) at a final concentration of 10 ng/ μ L and a total of 0.5 μ g of genomic DNA. DNA was quantified
696 using a Qubit double-stranded broad range DNA Assay or the high-sensitivity assay
697 (ThermoFisher Scientific) if necessary. Additionally, the integrity of the genomic DNA was
698 confirmed on agarose gels and the cleanness of the extracts tested by absence of inhibition during
699 PCR. For samples collected the following year, DNA was co-extracted with RNA (see next
700 section), in addition to extracting subsamples (10 g of soil) from the same core following the
701 extraction protocol described above ([Supplementary Table 1](#)).

702 Clean DNA extracts and co-extracts were submitted for sequencing at the Joint Genome
703 Institute (Walnut Creek, CA), where samples were subjected to a quality control check. Two of
704 the 96 samples from 2015 failed QC and thus were not sequenced (ERMZT233 and ERMZT446),
705 and four samples were sequenced ahead of the others (ERMLT700, ERMLT890, ERMZT100, and
706 ERMZT299). Ten out of 15 of the DNA co-extracts from 2016 failed QC due to low DNA yields
707 and were not sequenced either. Sequencing libraries for the first four samples were prepared in
708 microcentrifuge tubes. 100 ng of Genomic DNA was sheared to 600 bp pieces using the Covaris
709 LE220 and size selected with SPRI using AMPureXP beads (Beckman Coulter). The fragments
710 were treated with end-repair, A-tailing, and ligation of Illumina compatible adapters (IDT, Inc)
711 using the KAPA Illumina Library prep kit (KAPA biosystems). Libraries for the rest of the samples
712 were prepared in 96-well plates. Plate-based DNA library preparation for Illumina sequencing was
713 performed on the PerkinElmer Sciclone NGS robotic liquid handling system using Kapa
714 Biosystems library preparation kit. 200 ng of sample DNA was sheared to 600 bp using a Covaris
715 LE220 focused-ultrasonicator. The sheared DNA fragments were size selected by double-SPRI
716 and then the selected fragments were end-repaired, A-tailed, and ligated with Illumina compatible
717 sequencing adaptors from IDT containing a unique molecular index barcode for each sample
718 library.

719 All the libraries were quantified using KAPA Biosystem's next-generation sequencing
720 library qPCR kit and a Roche LightCycler 480 real-time PCR instrument. The quantified libraries

721 were then multiplexed with other libraries, and the pool of libraries was prepared for sequencing
722 on Illumina HiSeq sequencing platform utilizing a TruSeq paired-end cluster kit, v4, and
723 Illumina's cBot instrument to generate a clustered flow cell for sequencing. Sequencing of the
724 flow cell was performed on the Illumina HiSeq 2500 sequencer using HiSeq TruSeq SBS
725 sequencing kits, v4, following a 2x150 indexed run recipe.

726

727 *RNA-DNA co-extraction and sequencing*

728 Total RNA was extracted from a subset of 15 samples using the RNA PowerSoil Total
729 RNA isolation kit (Qiagen). Soil samples (8 g) preserved in LifeGuard solution (Qiagen) were
730 thawed on ice and centrifuged at 2,500 x g for 5 minutes to collect the soil at the bottom of the
731 tubes. As a supernatant, the LifeGuard solution was extracted from the tubes and aliquoted into
732 three 15 mL conical tubes that were used to transfer three separate 2 g subsamples for later use.
733 The remaining 2 g were split in half into two of the kit's bead tubes with pre-aliquoted bead
734 solution (to disperse the cells and soil particles). The lysis solution (SR1) and the non-DNA organic
735 and inorganic precipitation solution (SR2) were not added to the bead tube until all the subsamples
736 to be processed in a given day had been aliquoted. Subsamples were kept at - 20 °C before
737 transferring them to a - 80 °C freezer for permanent storage. The remainder of the extraction was
738 carried out following the manufacturer's instructions. An RNA PowerSoil DNA elution accessory
739 kit was used to co-extract DNA from the RNA capture columns, which was quantified as
740 previously described. A DNase treatment was performed in all the RNA extracts with a TURBO
741 DNA-free kit (Ambion) using 4 U of TURBO DNase at 37 °C for 30 minutes. The absence of
742 DNA was tested by PCR with universal primers to the SSU rRNA gene, and the integrity of the
743 RNA was checked using a Bioanalyzer RNA 6000 Nano kit following the manufacturer's
744 instructions. Total RNA was quantified before and after DNase treatments using a Qubit high-
745 sensitivity RNA assay (ThermoFisher Scientific). One of the RNA extracts (ERMLT590) did not
746 yield enough RNA for sequencing.

747 Total RNA and DNA co-extracts were submitted for sequencing at the Joint Genome
748 Institute in Walnut Creek, CA, where samples were subjected to a quality control check. rRNA
749 was removed from 1 µg of total RNA using Ribo-Zero(TM) rRNA Removal Kit (Illumina).
750 Stranded cDNA libraries were generated using the Illumina Truseq Stranded mRNA Library Prep
751 kit. The rRNA depleted RNA was fragmented and reversed transcribed using random hexamers

752 and SSII (Invitrogen) followed by second strand synthesis. The fragmented cDNA was treated
753 with end-pair, A-tailing, adapter ligation, and 8 cycles of PCR. For low input extracts, rRNA was
754 removed from 100 ng of total RNA using Ribo-Zero(TM) rRNA Removal Kit (Illumina). Stranded
755 cDNA libraries were generated using the Illumina Truseq Stranded mRNA Library Prep kit. The
756 rRNA depleted RNA was fragmented and reversed transcribed using random hexamers and SSII
757 (Invitrogen) followed by second strand synthesis. The fragmented cDNA was treated with end-
758 pair, A-tailing, adapter ligation, and 10 cycles of PCR. The prepared libraries were quantified using
759 KAPA Biosystem's next-generation sequencing library qPCR kit and run on a Roche LightCycler
760 480 real-time PCR instrument. The quantified libraries were then multiplexed with other libraries,
761 and the pool of libraries was prepared for sequencing on the Illumina HiSeq sequencing platform
762 utilizing a TruSeq paired-end cluster kit, v4, and Illumina's cBot instrument to generate a clustered
763 flow cell for sequencing. Sequencing of the flow cell was performed on the Illumina HiSeq 2500
764 sequencer using HiSeq TruSeq SBS sequencing kits, v4, following a 2 x 150 indexed run recipe.

765

766 *Metagenomes assembly and annotation and ribosomal protein L6 analysis*

767 Methods used for 2015 and 2016 metagenomes assembly and annotation are described
768 elsewhere ²⁷. In brief, after quality filtering, reads from individual samples were assembled
769 separately using IDBA-UD v1.1.1 with a minimum k-mer size of 40, a maximum k-mer size of
770 140 and step size of 20. Only contigs > 1Kb were kept for further analyses. Gene prediction was
771 done with Prodigal v2.6.3 in meta mode, annotations obtained using USEARCH against Uniprot,
772 Uniref90 and KEGG, and 16S rRNA and tRNAs predicted as described in Diamond et al. ⁵. Reads
773 were mapped to the assemblies using Bowtie2 ²⁸ and default settings to estimate coverage. To
774 estimate the number of genomes potentially present across all 94 metagenomes, we used the
775 ribosomal protein L6 as marker gene and RPxSuite
776 (<https://github.com/alexcritschroph/RPxSuite>) as described in Olm et al. ⁶.

777

778 *Genome binning, curation, and dereplication*

779 Annotated metagenomes from both years were uploaded onto ggkbase
780 (<https://ggkbase.berkeley.edu>), where binning tools based on GC content, coverage and winning
781 taxonomy ²⁹ were used for genome binning. These bins and additional bins that were obtained with
782 the automated bidders ABAWACA1 (<https://github.com/CK7/abawaca>), ABAWACA2,

783 MetaBAT ³⁰, Maxbin2 ³¹ and Concoct ³² were pooled, and DASStool was used for selection of the
784 best set of bins from each sample as described by Diamond et al. ⁵. Notably, no bins were recovered
785 from sample ERMZT266 by any method.

786 Genomic bins were filtered based on completeness $\geq 70\%$ of a set of 51 bacterial single
787 copy genes (BSCG) if affiliated with Bacteria and a set of 38 archaeal single copy genes (ASCG);
788 and a level of contamination $\leq 10\%$ based on the corresponding list of single copy genes ³³.
789 Additionally, bins that were 59-68% complete with a highest taxonomic level defined as Bacteria
790 in ggKbase, or potential members of the candidate phyla radiation (CPR) were kept for further
791 scrutiny. To obtain a set of genomes for visual curation in ggKbase, genomes were dereplicated at
792 99% ANI across samples located within a given floodplain using dRep with the --
793 ignoreGenomeQuality flag. ³⁴. Any assembly error in the dereplicated set was addressed using
794 ra2.py ³⁵, and contigs that fell below the 1 Kb length minimum after this step were removed from
795 the bins. At this point, the level of completeness of CPR genomes was confirmed based on a list
796 of 43 BSCG ⁷. Genomes that did not meet the completeness thresholds post-assembly error
797 correction and that were not affiliated with CPR or novel bacteria were removed from the analysis.
798 Considering that bins changed as a result of this process, genes were re-predicted using Prodigal
799 in single mode, reads were mapped to the bins using Bowtie2, and bins were re-imported onto
800 ggKbase. Visual inspection of taxonomic profile, GC content and to a minor extent coverage,
801 allowed us to further reduce contamination. The final set of 248 curated bins from 2015 was
802 dereplicated at 98% ANI this time across floodplains including the --genomeInfo flag to take into
803 account completeness and contamination in the process of representative bin selection. Within this
804 set, genomes $\geq 90\%$ complete were deemed near-complete ([Supplementary Table 2](#)). Eight
805 relatively low coverage genomes fell just below the completeness requirement due to
806 fragmentation after curation to remove possible local assembly errors; these were retained as they
807 represent important taxonomic diversity.

808 Similarly, genomes reconstructed from floodplain L samples collected in 2016 that passed
809 the completeness ($\geq 70\%$) and contamination thresholds ($\leq 10\%$) were visually inspected and
810 improved in ggKbase. Assembly errors were corrected with ra2.py ³⁵, and contigs that fell below
811 the 1Kb length were removed, as well as genomes that did not pass the thresholds for completeness
812 after assembly error correction. Genes were re-predicted using Prodigal in single mode and the
813 final set of curated genomes were imported onto ggKbase.

814 To determine whether the same species were present in two different years, we pooled the
815 genome set from 2015 and the curated 2016 set and dereplicated using dRep at 95% ANI including
816 the --genomeInfo flag to take into account completeness and contamination in the process of
817 representative bin selection³⁴. In this set of genomes, 13 were reconstructed from a deeper depth
818 ([Supplementary Table 3](#)). However, only 3 genomes were unique and the other 10 clustered with
819 genomes reconstructed from the ~10-25 cm depth, indicating overlap between the species found
820 at the two depths. Therefore, we kept these genomes for further analyses.

821

822 *Genome metabolic annotation*

823 We carefully chose a set of ecologically relevant proteins that catalyze geochemical
824 transformations related to aerobic respiration, metabolism of C1 compounds, hydrogen
825 metabolism, nitrogen cycling, and sulfur cycling ([Supplementary Table 4](#)). Hidden Markov
826 Models (HMMs) for the majority of these proteins were obtained from KOfam, the customized
827 HMM database of KEGG Orthologs (KOs)³⁶. Custom-made HMMs targeting nitrite
828 oxidoreductase subunits A and B (NxrA and NxrB), periplasmic cytochrome *c* nitrite reductase
829 (NirS, cd1-NIR heme-containing), cytochrome *c*-dependent nitric oxide reductase (NorC; cNOR),
830 hydrazine dehydrogenase (HzoA), hydrazine synthase (HzsA), dissimilatory sulfite reductase D
831 (DsrD), sulfide:quinone reductase (Sqr), sulfur dioxygenase (Sdo), ribulose-bisphosphate
832 carboxylase (RuBisCO) form I and form II, and alcohol dehydrogenases (Pqq-XoxF-MxaF) were
833 obtained from Anantharaman et al.⁷. NiFe and FeFe hydrogenases were predicted using HMMs
834 from Méheust et al.³⁷ and assigned to functional groups following Matheus Carnevali et al.²⁹ (see
835 Phylogenetic Analyses below for tree construction methods; [Supplementary Data 3 and 4 and](#)
836 [Supplementary Tables 9 and 10](#)). No real group 4 membrane bound NiFe hydrogenases were
837 identified among the East River representative genomes (data not shown). HMMER3³⁸ was used
838 to annotate the dereplicated sets of genomes following predefined score cutoffs³⁶. A subset (10%)
839 of the hits to all of these HMMs were visually checked to determine whether the cutoffs were
840 appropriate for this dataset as described in Lavy et al.³⁹ and Jaffe et al.⁴⁰. Only in the case of
841 formate dehydrogenase (FdhA (K05299 and K22516), FdoG/FdhF/FdwA (K00123)) the cutoff
842 was lowered to include additional hits.

843 For a protein to be considered potentially encoded in the genome, the catalytic subunit and
844 the majority of the accessory subunits had to be detected by the corresponding HMMs at the

845 established cutoffs. The implication for these function definitions is that in some cases even if
846 some subunits that make up an enzyme were detected, the enzyme could have been deemed absent
847 because a key part was missing (Supplementary Table 4). Similarly, pathways that require the
848 activity of multiple enzymes were only detectable if all of the enzymes were present. Only in cases
849 like the Wood-Jungdahl pathway we required the majority of the genes to be present, taking into
850 consideration genome completeness. Furthermore, if multiple enzymes could catalyze a given
851 reaction (e.g., use O₂ as a terminal electron acceptor) the presence of genes encoding one such
852 enzyme in a genome would be indicative that this capacity was present in the genome.
853 Additionally, if different pathways lead to the same biogeochemical transformation (e.g., CO₂-
854 fixation), the presence of genes encoding one of those pathways (or key enzymes) was considered
855 as sufficient to indicate its presence (Supplementary Table 4). In a limited number of cases a given
856 pathway may also involve enzymes that are part of central metabolism or that are part of multiple
857 pathways, and in these cases we chose to define presence based on the key catalyst instead of the
858 whole pathway (e.g., RuBisCO in the Calvin Benson pathway).

859 Carbohydrate active enzymes were predicted using the Carbohydrate-Active enZymes
860 Database (CAZY; <http://www.cazy.org/>)¹⁰ (version 1.0) (e-value cut-off 1e-20).

861

862 *Genome coverage and detection*

863 Reads were mapped to the dereplicated set of bins using Bowtie2²⁸ and a mismatch
864 threshold of 2% dissimilarity. Calculate_coverage.py
865 (<https://github.com/christophertbrown/bioscripts/tree/master/ctbBio>) was used to estimate the
866 average number of reads mapping to each genome and the proportion of the genome that was
867 covered by reads (breadth). Genomes with a coverage of at least 0.01 X were considered to be
868 detected in a given sample. The Hellinger transformation was used to account for differences in
869 sequencing depth among samples and determine final genome abundance. To illustrate genome
870 detection across samples we used the ggplot2 package⁴¹. Genomes were clustered by average
871 linkage using the Hellinger transformed abundance across samples (from read mapping), and the
872 samples were clustered by Euclidean distance in R⁴².

873

874 *Phylogenetic analyses*

875 Two phylogenetic trees were constructed with a set of 14 ribosomal proteins (L2, L3, L4,
876 L5, L6, L14, L15, L18, L22, L24, S3, S8, S17, and S19). One tree included Betaproteobacteria
877 genomes from this study at the subspecies level (98% ANI) and ~ 1540 reference
878 Betaproteobacteria genomes from the NCBI ([Supplementary Figure 2 and Supplementary Data 1](#)).
879 The other tree included the set of 215 genomes dereplicated at 95% ANI and ~ 2,228 reference
880 genomes from the NCBI genome database ([Supplementary Data 2](#)). For each genome, the
881 ribosomal proteins were collected along the scaffold with the highest number of ribosomal
882 proteins. A maximum-likelihood tree was calculated based on the concatenation of the ribosomal
883 proteins as follows: Homologous protein sequences were aligned using MAFFT (version 7.390) (-
884 -auto option)⁴³, and alignments refined to remove gapped regions using Trimal (version 1.4.22) (-
885 -gappyout option)⁴⁴. Tree reconstruction was performed using IQ-TREE (version 1.6.12) (as
886 implemented on the CIPRES web server⁴⁵, using ModelFinder⁴⁶ to select the best model of
887 evolution (LG+I+G4), and with 1000 ultrafast bootstrap⁴⁷. Taxonomic affiliations were
888 determined based on the closest reference sequences relative to the query sequences on the tree
889 and extended to other members of the ANI cluster. In many cases, the phylogeny was not clear
890 upon first inspection of the tree and additional reference genomes were added if publicly available.
891

892 Phylogenetic trees for proteins of interest were reconstructed using the same methods
893 described above, except with different sets of reference sequences. East River homologs in the
894 dimethyl sulfoxide reductase (DMSOR) superfamily such as the catalytic subunit of formate
895 dehydrogenase (FdhA), nitrite oxidoreductase (NxrA), membrane-bound nitrate reductase (NarG;
896 H⁺-translocating), and periplasmic nitrate reductase subunit A (NapA) were confirmed by
897 phylogeny on a tree with reference sequences from Méheust et al.³⁷ ([Supplementary Table 11 and](#)
898 [Supplementary Data 5](#)). To distinguish form I and form II CODHs and other other subtypes among
899 homologs to K03520 we used Diamond's et al.⁵ dataset, which comprises reference sequences
900 from Quiza et al.¹⁶ ([Supplementary Table 12 and Supplementary Data 6](#)). Similarly, homologs
901 identified using the Pqq-XoxF-MxaF HMM for alcohol dehydrogenases were placed on a
902 phylogenetic tree with reference sequences from Diamond's et al.⁵ dataset, comprising references
903 from Keltjens et al.⁴⁸ and Taubert et al.⁴⁹. In this tree, all East River homologs were clustered
904 with methanol dehydrogenases ([Supplementary Table 13 and Supplementary Data 7](#)) instead of
905 other types of alcohol dehydrogenases. To distinguish between dissimilatory (bi)sulfite reductase

906 oxidative or reductive bacterial types, DsrA and DsrB homologs from individual genomes were
907 concatenated to each other, aligned, and added to a phylogenetic tree with reference sequences
908 from Muller et al.⁵⁰ ([Supplementary Table 14 and Supplementary Data 8](#)).

909

910 *Community diversity and composition*

911 Diversity indices for each sample were calculated from the Hellinger transformed
912 abundance table for the genome set at subspecies level (98% ANI) using the vegan package in R
913 ⁵¹. Species numbers and Shannon diversity per sample were quantified using the specnumber and
914 vegdist functions of vegan respectively ([Supplementary Figure 3](#)). An analysis of variance,
915 implemented in the aov function in R, was used to test for significant differences in mean species
916 number and Shannon diversity in relationship to the floodplain samples originated from. No
917 significant differences in group means were detected considering a p-value < 0.05 as significant.

918 To investigate community composition at the phylum/class level as determined by
919 phylogenetic analysis, the Hellinger-transformed abundance table for the genome set at the
920 subspecies level (98% ANI) was converted to a presence/absence table. The number of samples
921 where each genome was detected was counted and the number of genomes affiliated to a given
922 taxon was summed by sample and plotted in R with ggplot2 ⁴¹.

923

924 *Identification of a core floodplain microbiome*

925 To identify organisms that were a “core” or “shared” set across all sampled sites, we
926 operationally defined a core set as: (1) organisms that were not statistically associated with any
927 specific floodplain using indicator species analysis, and (2) who were detected (displayed $\geq 0.01X$
928 coverage) in at least 89 of the 94 total samples (the 90th percentile for this level of presence across
929 all 248 genomes). Indicator species analysis was performed on the log transformed coverage
930 values that were filtered to include only coverage values $\geq 0.01X$ using the indicpecies package
931 ⁵² in R version 3.5.2 (R core team 2018) ⁴² with 9999 permutations. All p-values for associations
932 of an organism genome with a floodplain or group of floodplains were then subsequently corrected
933 using False Discovery Rate with $FDR \leq 0.05$ being considered a significant association. This
934 resulted in 42 genomes that were not statistically associated with any floodplain by ISA and were
935 also detected in ≥ 89 samples ([Supplementary Table 5](#)). For visualization of organism abundance
936 profiles in relationship to their membership in the core floodplain microbiome, ISA clusters, and

937 relative to the coefficient of variation of their coverage, Hellinger normalized coverage data was
938 projected onto a two dimensional space using Uniform Manifold Approximation and Projection
939 (UMAP) implemented in the uwot package in R ⁵³ using the following parameters: `umap(data =`
940 `coverage_data, n_neighbors = 15, nn_method = "fnn", spread = 5, min_dist = 0.01, n_components`
941 `= 2, metric = "euclidian", n_epochs = 1000)`.

942

943 *Identification of enriched metabolic functions in core floodplain microbiome*

944 Overrepresentation of metabolic functions within the set of genomes comprising the core
945 floodplain microbiome (n = 42) was assessed using hypergeometric testing. The probability of
946 observing the number of genomes in the core floodplain microbiome carrying each of 33 functions,
947 given the total number of genomes with that function across our full genomic dataset (n = 248),
948 was calculated using the `phyper` function in R. Probabilities calculated across all metabolic
949 functions were corrected for multiple testing using false discovery rate with the `p.adjust` function
950 in R and with $FDR \leq 0.05$ being considered a significant enrichment of a function in the core
951 microbiome.

952

953 *Analysis of correlations among environmental variables*

954 Correlations between numeric soil biogeochemical variables across samples were
955 calculated using spearman rank correlation implemented in the `rcorr` function of the `Hmisc`
956 package in R (<https://github.com/harrelfe/Hmisc>). Correlations between variables were then plotted as
957 a correlogram and ordered using hierarchical clustering with Ward's method using the `corrplot`
958 package in R ⁵⁴.

959

960 *Fourth corner analysis*

961 A `rlq`-fourth corner analysis was performed on genome abundances, environmental data,
962 and genome metabolic annotations using the R package `ade4` ⁵⁵. Specifically, the pre-Hellinger
963 transformed genome abundance table was used for a correspondence analysis, the selected
964 environmental variables (see *Soil Chemistry* and *GIS*) were used for a Hill-Smith analysis, and the
965 genome metabolic annotations were used for PCA. A randomization test (as described by ter Braak
966 et al. ⁵⁶ and Dray et al. ⁵⁷) was used to test the global significance of the trait-environment
967 relationships. The fourth-corner statistic was then calculated on the same inputs as the `rlq` analysis

968 with 50,000 permutations and p-value adjustments using the FDR global methods. The results of
969 the rlq-fourth corner analysis were plotted using the ggplot2 package ⁴¹.

970

971 *Metatranscriptomic analyses*

972 To determine differentially transcribed genes, potential levels of activity by phylum or
973 class, most transcribed CAZY, and most transcribed genes among key geochemical
974 transformations, metatranscriptomic reads were mapped using Bowtie2 ²⁸ to a set of high-quality
975 draft genomes dereplicated at 95% (see above). Read pairs were then filtered by a minimum
976 identity of 95% to the reference with MAPQ \geq 2 and total number of mapped read pairs was
977 counted for each gene. Counts for metabolic genes were analyzed with DESeq2 ⁹ to determine
978 differential expression in response to soil organic carbon and p-values were adjusted to correct for
979 multiple hypothesis testing (FDR $<$ 0.05).

980

981 *GIS*

982 All GIS operations and cartographic visualizations were performed in QGIS v2.12.1 except
983 where otherwise stated. The base remote sensed imagery used was obtained from USDA NAIP
984 (USDA-FSA Aerial Photography Field Office publication date 20171220; 1m ground pixel
985 resolution). Digital terrain model (DTM) at a ground resolution of 0.5 m/pixel was derived by
986 airborne LiDAR data acquired by Quantum Spatial in collaboration with Eagle Mapping Ltd ⁵⁸
987 (doi:10.21952/WTR/1412542) in 2015. All maps were projected using EPSG:26913 NAD83/
988 UTM zone 13N. Meander and adjacent river polygons were manually delineated in QGIS. The
989 distance from a sample point to the manually delineated river polygons was calculated using the
990 NNJoin tool. To calculate the sample distances to meander toe, lines were manually drawn
991 between all samples and the meander toe perpendicular to river flow and distances calculated using
992 NNJoin ([Supplementary Figure 5](#)). Similarly, to calculate sample distances to the middle of the
993 meander, a line perpendicular to the meander toe line was drawn across the middle of the meander
994 ([Supplementary Figure 5](#)). Sample distances to this line were also calculated using NNJoin and
995 samples on the downstream side of the line were converted to negative values to indicate upstream
996 and downstream sides of the meander. TPI is computed from the DTM as the difference between
997 the elevation of a center point and the average elevation measured in the neighboring area (3 by 3
998 m) ⁵⁹. To display genome abundances as used in the rlq-fourth corner analysis, filtered abundance

999 values were chi-square transformed in R using the *decostand* in the *vegan* package and exported
1000 to display in QGIS. Spatial kriging of inorganic carbon was performed in R. The manually
1001 delineated meander polygons were converted to *SpatialPixelsDataFrame* using the *sp* package. A
1002 simple variogram model was fit to the natural log transformed inorganic carbon values with a
1003 spatial cutoff of 60 m. Kriging was then performed using the sample points, the meander
1004 *SpatialPixelsDataFrame*, and the fitted variogram model. The natural log transformed inorganic
1005 carbon values were then back transformed and the kriged map exported for visualization in QGIS.

1006

1007 **Data Availability**

1008 Representative genomes in the subspecies level set (98% ANI) can be accessed at
1009 https://ggkbase.berkeley.edu/ER15_ALL_curated_dRep98/organisms and representative
1010 genomes in the species level set (95% ANI) can be accessed at
1011 https://ggkbase.berkeley.edu/ER15ALL_ERML16_dRep95/organisms. Please note ggKbase is a
1012 'live' site, genomes may be updated after this publication. Raw sequence reads for all metagenomes
1013 and metatranscriptomes included in this study can be accessed in the NCBI Bioproject Database
1014 using the umbrella accession number (PRJNA630765). Supplementary Table 1 includes NCBI
1015 Bioproject accession numbers for individual metagenomes metatranscriptomes.

1016

1017 **Acknowledgements**

1018 We are grateful to Chad Hobson, David McGrath, and Rosemary Carrol for support in the
1019 field during samples collection; and to Joel Rowland (Los Alamos National Lab) and Helen
1020 Malenda (USGS) for useful insights. We thank the Rocky Mountain Biological Laboratory for lab
1021 space at the field site. This work was supported as part of the Watershed Function Scientific Focus
1022 Area funded by the U.S. Department of Energy, Office of Science, Office of Biological and
1023 Environmental Research under Award Number DE-AC02-05CH11231. Sequencing was
1024 conducted at the Joint Genome Institute (a DOE Office of Science User Facility) under a CSP
1025 award.

1026

1027 **References**

1028

- 1029 1. Viviroli D, Dürr HH, Messerli B, Meybeck M, Weingartner R. Mountains of the world, water
1030 towers for humanity: Typology, mapping, and global significance. *Water Resour Res* **43**, (2007).

1031

- 1032 2. Immerzeel WW, *et al.* Importance and vulnerability of the world's water towers. *Nature* **577**, 364-
1033 369 (2020).
1034
- 1035 3. Hubbard SS, *et al.* The East River, Colorado, Watershed: A Mountainous Community Testbed for
1036 Improving Predictive Understanding of Multiscale Hydrological–Biogeochemical Dynamics.
1037 *Vadose Zone J* **17**, 180061 (2018).
1038
- 1039 4. Levin SA. The Problem of Pattern and Scale in Ecology: The Robert H. MacArthur Award Lecture.
1040 *Ecology* **73**, 1943-1967 (1992).
1041
- 1042 5. Diamond S, *et al.* Mediterranean grassland soil C–N compound turnover is dependent on rainfall
1043 and depth, and is mediated by genomically divergent microorganisms. *Nat Microbiol* **4**, (2019).
1044
- 1045 6. Olm MR, Crits-Christoph A, Diamond S, Lavy A, Matheus Carnevali PB, Banfield JF. Consistent
1046 Metagenome-Derived Metrics Verify and Delineate Bacterial Species Boundaries. *mSystems* **5**,
1047 e00731-00719 (2020).
1048
- 1049 7. Anantharaman K, *et al.* Thousands of microbial genomes shed light on interconnected
1050 biogeochemical processes in an aquifer system. *Nat Commun* **7**, 13219 (2016).
1051
- 1052 8. Dray S, *et al.* Combining the fourth-corner and the RLQ methods for assessing trait responses to
1053 environmental variation. *Ecology* **95**, 14-21 (2014).
1054
- 1055 9. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq
1056 data with DESeq2. *Genome Biol* **15**, 550 (2014).
1057
- 1058 10. Lombard V, Golaconda Ramulu H, Drula E, Coutinho PM, Henrissat B. The carbohydrate-active
1059 enzymes database (CAZy) in 2013. *Nucleic Acids Res* **42**, D490-D495 (2014).
1060
- 1061 11. Dwivedi D, *et al.* Geochemical Exports to River From the Intrameander Hyporheic Zone Under
1062 Transient Hydrologic Conditions: East River Mountainous Watershed, Colorado. *Water Resour*
1063 *Res* **54**, 8456-8477 (2018).
1064
- 1065 12. Fox PM, *et al.* Shale as a Source of Organic Carbon in Floodplain Sediments of a Mountainous
1066 Watershed. *J Geophys Res: Biogeosci* **125**, e2019JG005419 (2020).
1067
- 1068 13. Argiroff WA, Zak DR, Lanser CM, Wiley MJ. Microbial Community Functional Potential and
1069 Composition Are Shaped by Hydrologic Connectivity in Riverine Floodplain Soils. *Microb Ecol*
1070 **73**, 630-644 (2017).
1071
- 1072 14. Wan J, *et al.* Predicting sedimentary bedrock subsurface weathering fronts and weathering rates.
1073 *Sci Rep* **9**, 17198 (2019).
1074
- 1075 15. King GM, Weber CF. Distribution, diversity and ecology of aerobic CO-oxidizing bacteria. *Nat*
1076 *Rev Microbiol* **5**, 107-118 (2007).
1077
- 1078 16. Quiza L, Lalonde I, Guertin C, Constant P. Land-use influences the distribution and activity of high
1079 affinity CO-oxidizing bacteria associated to type I-coxL genotype in soil. *Front Microbiol* **5**,
1080 (2014).
1081

- 1082 17. Cordero PRF, *et al.* Atmospheric carbon monoxide oxidation is a widespread mechanism
1083 supporting microbial survival. *ISME J* **13**, 2868-2881 (2019).
1084
- 1085 18. Boutard M, *et al.* Functional Diversity of Carbohydrate-Active Enzymes Enabling a Bacterium to
1086 Ferment Plant Biomass. *PLoS Genet* **10**, e1004773 (2014).
1087
- 1088 19. Houlton BZ, Morford SL, Dahlgren RA. Convergent evidence for widespread rock nitrogen sources
1089 in Earth's surface environment. *Science* **360**, 58-62 (2018).
1090
- 1091 20. Winnick MJ, Carroll RWH, Williams KH, Maxwell RM, Dong W, Maher K. Snowmelt controls
1092 on concentration-discharge relationships and the balance of oxidative and acid-base weathering
1093 fluxes in an alpine catchment, East River, Colorado. *Water Resour Res* **53**, 2507-2523 (2017).
1094
- 1095 21. Bernal S, Hedin LO, Likens GE, Gerber S, Buso DC. Complex response of the forest nitrogen cycle
1096 to climate change. *Proc Natl Acad Sci USA* **109**, 3406-3411 (2012).
1097
- 1098 22. Sebestyen SD, *et al.* Sources, transformations, and hydrological processes that control stream
1099 nitrate and dissolved organic matter concentrations during snowmelt in an upland forest. *Water*
1100 *Resour Res* **44**, (2008).
1101
- 1102 23. Rütting T, Boeckx P, Müller C, Klemmedtsson L. Assessment of the importance of dissimilatory
1103 nitrate reduction to ammonium for the terrestrial nitrogen cycle. *Biogeosciences* **8**, 1779-1791
1104 (2011).
1105
- 1106 24. Hefting M, *et al.* Water table elevation controls on soil nitrogen cycling in riparian wetlands along
1107 a European climatic gradient. *Biogeochemistry* **67**, 113-134 (2004).
1108
- 1109 25. Pribulick CE, *et al.* Contrasting the hydrologic response due to land cover and climate change in a
1110 mountain headwaters system. *Ecohydrology* **9**, 1431-1438 (2016).
1111
- 1112 26. Albertsen M, Hugenholtz P, Skarshewski A, Nielsen KL, Tyson GW, Nielsen PH. Genome
1113 sequences of rare, uncultured bacteria obtained by differential coverage binning of multiple
1114 metagenomes. *Nat Biotechnol* **31**, 533-538 (2013).
1115
- 1116 27. Lavy A, *et al.* Microbial communities across a hillslope-riparian transect shaped by proximity to
1117 the stream, groundwater table, and weathered bedrock. *Ecol Evol* **0**, (2019).
1118
- 1119 28. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**, 357-359
1120 (2012).
1121
- 1122 29. Matheus Carnevali PB, *et al.* Hydrogen-based metabolism as an ancestral trait in lineages sibling
1123 to the Cyanobacteria. *Nat Commun* **10**, 463 (2019).
1124
- 1125 30. Kang DD, Froula J, Egan R, Wang Z. MetaBAT, an efficient tool for accurately reconstructing
1126 single genomes from complex microbial communities. *PeerJ* **3**, e1165 (2015).
1127
- 1128 31. Wu Y-W, Simmons BA, Singer SW. MaxBin 2.0: an automated binning algorithm to recover
1129 genomes from multiple metagenomic datasets. *Bioinformatics* **32**, 605-607 (2016).
1130
- 1131 32. Alneberg J, *et al.* Binning metagenomic contigs by coverage and composition. *Nat Methods* **11**,
1132 1144-1146 (2014).

- 1133
1134 33. Sieber CMK, *et al.* Recovery of genomes from metagenomes via a dereplication, aggregation and
1135 scoring strategy. *Nat Microbiol* **3**, 836-843 (2018).
1136
1137 34. Olm MR, Brown CT, Brooks B, Banfield JF. dRep: a tool for fast and accurate genomic
1138 comparisons that enables improved genome recovery from metagenomes through de-replication.
1139 *ISME J* **11**, 2864-2868 (2017).
1140
1141 35. Brown CT, *et al.* Unusual biology across a group comprising more than 15% of domain Bacteria.
1142 *Nature* **523**, 208-211 (2015).
1143
1144 36. Aramaki T, *et al.* KofamKOALA: KEGG ortholog assignment based on profile HMM and adaptive
1145 score threshold. *bioRxiv*, 602110 (2019).
1146
1147 37. Méheust R, *et al.* Aquatic Elusimicrobia are metabolically diverse compared to gut microbiome
1148 Elusimicrobia and some have novel nitrogenase-like gene clusters. *bioRxiv*, 765248 (2019).
1149
1150 38. Pagnuco IA, Revuelta MV, Bondino HG, Brun M, ten Have A. HMMER Cut-off Threshold Tool
1151 (HMMERCTTER): Supervised classification of superfamily protein sequences with a reliable cut-
1152 off threshold. *PloS one* **13**, e0193757 (2018).
1153
1154 39. Lavy A, *et al.* Taxonomically and metabolically distinct microbial communities with depth and
1155 across a hillslope to riparian zone transect. *bioRxiv*, 768572 (2019).
1156
1157 40. Jaffe AL, Castelle CJ, Carnevali PBM, Gribaldo S, Banfield JF. The rise of diversity in metabolic
1158 platforms across the Candidate Phyla Radiation. *bioRxiv*, (2019).
1159
1160 41. Wickham H. *ggplot2: elegant graphics for data analysis*. Springer (2016).
1161
1162 42. Team RC. R: a language and environment for statistical computing. R Foundation for Statistical
1163 Computing. 2014. (2015).
1164
1165 43. Katoh K, Standley DM. A simple method to control over-alignment in the MAFFT multiple
1166 sequence alignment program. *Bioinformatics* **32**, 1933-1942 (2016).
1167
1168 44. Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. trimAl: a tool for automated alignment
1169 trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972-1973 (2009).
1170
1171 45. Miller MA, Pfeiffer W, Schwartz T. *Creating the CIPRES Science Gateway for inference of large*
1172 *phylogenetic trees*. In: *Proceedings of the Gateway Computing Environments Workshop (GCE)*
1173 (2010).
1174
1175 46. Kalyaanamoorthy S, Minh BQ, Wong TK, von Haeseler A, Jermini LS. ModelFinder: fast model
1176 selection for accurate phylogenetic estimates. *Nat Methods* **14**, 587 (2017).
1177
1178 47. Hoang DT, Chernomor O, Von Haeseler A, Minh BQ, Vinh LS. UFBoot2: improving the ultrafast
1179 bootstrap approximation. *Mol Biol Evol* **35**, 518-522 (2018).
1180
1181 48. Keltjens JT, Pol A, Reimann J, Op den Camp HJM. PQQ-dependent methanol dehydrogenases:
1182 rare-earth elements make a difference. *Appl Microbiol Biotechnol* **98**, 6163-6183 (2014).
1183

- 1184 49. Taubert M, *et al.* XoxF encoding an alternative methanol dehydrogenase is widespread in coastal
1185 marine environments. *Environ Microbiol* **17**, 3937-3948 (2015).
1186
- 1187 50. Müller AL, Kjeldsen KU, Rattei T, Pester M, Loy A. Phylogenetic and environmental diversity of
1188 DsrAB-type dissimilatory (bi)sulfite reductases. *ISME J* **9**, 1152-1165 (2014).
1189
- 1190 51. Oksanen J, *et al.* Vegan: Community ecology package. R package version 2.2-0 (pp. 10). (2014).
1191
- 1192 52. Cáceres MD, Legendre P. Associations between species and groups of sites: indices and statistical
1193 inference. *Ecology* **90**, 3566-3574 (2009).
1194
- 1195 53. Melville J. uwot: The uniform manifold approximation and projection (UMAP) method for
1196 dimensionality reduction. *R package version 00 09010*, (2019).
1197
- 1198 54. Wei T, Simko V. R package" corrplot": Visualization of a Correlation Matrix (Version 0.84).
1199 *Retrieved from, <https://github.com/taiyun/corrplot>*, (2017).
1200
- 1201 55. Dray S, Dufour A-B. The ade4 package: implementing the duality diagram for ecologists. *J Stat*
1202 *Softw* **22**, 1-20 (2007).
1203
- 1204 56. Ter Braak CJ, Cormont A, Dray S. Improved testing of species traits–environment relationships in
1205 the fourth-corner problem. *Ecology* **93**, 1525-1526 (2012).
1206
- 1207 57. Dray S, *et al.* Combining the fourth-corner and the RLQ methods for assessing trait responses to
1208 environmental variation. *Ecology* **95**, 14-21 (2014).
1209
- 1210 58. Wainwright H, & Williams, Kenneth. . LiDAR collection in August 2015 over the East River
1211 Watershed, Colorado, USA. United States. (doi:10.21952/WTR/1412542.) (2017).
1212
- 1213 59. Falco N, *et al.* Integrated imaging of above and below ground properties and their interactions: A
1214 case study in East River Watershed, Colorado. In: *SEG Technical Program Expanded Abstracts*
1215 *2018* (2018).
1216
1217
1218